

**IN THE UNITED STATES DISTRICT COURT
FOR THE EASTERN DISTRICT OF VIRGINIA**

THROUGHPUTER, INC.,

Plaintiff,

v.

MICROSOFT CORPORATION,

Defendant.

Civil Action No. 3:21cv216

JURY TRIAL DEMANDED

COMPLAINT FOR PATENT INFRINGEMENT

ThroughPuter, Inc. (“ThroughPuter”) hereby alleges for its Complaint (“Complaint”) against Microsoft Corp. (“Microsoft”) as follows.

INTRODUCTION

1. A Field Programmable Gate Array (FPGA) is a specific type of microprocessor that can be reconfigured based on the tasks to be performed by it. In some situations, this reconfiguration takes place in between the performance of tasks by the FPGA. For example, the FPGAs can be configured and re-configured to provide hardware accelerators for data processing functions such as encoding/decoding, encryption/decryption or compression/decompression and then reconfigured to perform speech translation in addition to, *e.g.*, encryption and compression. Starting in 2015, Microsoft started putting FPGA processors in each server for Microsoft Azure (“Azure”), which is the world’s largest cloud computing platform-as-a-service (PaaS). Using these FPGA processors for the underlying dynamic parallel execution environment or architecture, Microsoft claims a 150-200 fold improvement in data throughput and a 50 fold improvement in energy efficiency. Exhibit 15 at 91, Exhibit 29. Latency has also lowered by about a factor of 10. Exhibit 15 at 86, Exhibit 29. The end result is a power savings in Microsoft and increased processing speeds for Microsoft and the users of applications running on the world’s largest cloud computing platform.

2. In 2013, Plaintiff ThroughPuter disclosed a reconfigurable and dynamic parallel execution architecture running on FPGA processors in writing to Microsoft. Two years later, Microsoft filed a patent application on the same hardware-based fabric (failing to disclose to the Patent Office any information about ThroughPuter’s technology or earlier patent filings). Microsoft was ultimately awarded the patent, which has claims that match almost exactly those of ThroughPuter’s U.S. Patent No. 10,963,306, which is the basis for Count 1 of this complaint. In other words, Microsoft was awarded a patent on the same hardware-based fabric claimed in the

ThroughPuter patents, wrongly suggesting that Microsoft had invented the architecture underlying Azure where ThroughPuter had already patented that architecture and disclosed it to Microsoft.

3. The figure below left is an excerpt from the materials disclosed by ThroughPuter to Microsoft in 2013, disclosing an application load and type adaptive “manycore fabric.” The different colors show different reconfigurable cores (*e.g.*, logic blocks of FPGA processors) being assigned different tasks over time. The figure on the right is a colorized version of figures from Microsoft’s 2015 patent filing. The Microsoft figures show the same reconfiguration of cores to different tasks over time.

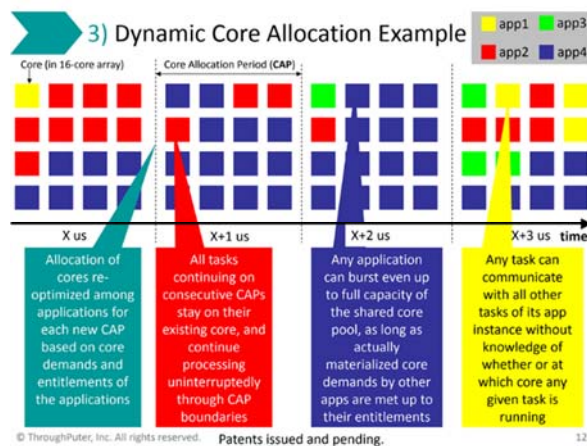


Exhibit 23 at 18.

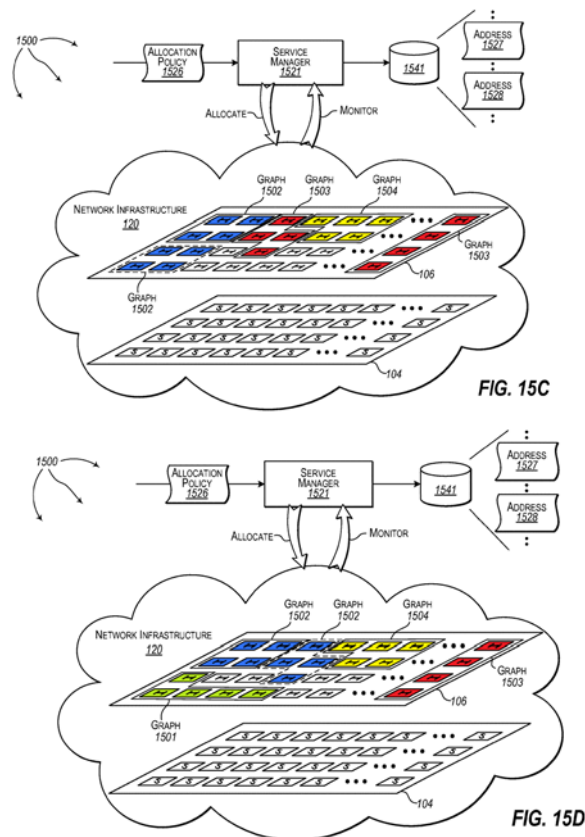


Exhibit 11, FIGs. 15C and 15D (annotated)

4. Microsoft’s copying of ThroughPuter’s technology was deliberate and wanton. While ThroughPuter has managed to survive as a small company developing and offering for sale other innovative solutions enabled by the throughput improvements which result from

ThroughPuter's patented technology, Microsoft's copying and infringement of ThroughPuter's intellectual property has irreparably damaged ThroughPuter's prospects of entering the PaaS market, which had grown from \$3.8B in 2015 to \$37.5B in 2019.¹

5. ThroughPuter has attempted repeatedly to engage with Microsoft to resolve this matter outside of litigation. Microsoft has rejected or ignored ThroughPuter's proposals and requests for business resolution of the matter.

6. Microsoft's decision to reject ThroughPuter's proposals and offers to collaborate has created a highly inequitable situation by which Microsoft has used its tremendous market power to scale up the world's largest and most successful cloud computing platform. Microsoft did this based on the technology it copied from ThroughPuter: a small start-up that could have only competed based on its effort to protect and patent its innovations. Microsoft vaporized the competitive edge to which ThroughPuter was entitled by its decision to ignore ThroughPuter while generating billions of dollars in revenue per quarter based on its unlawful exploitation of ThroughPuter's technology.

7. In a 2016 presentation entitled "Catapult at Ignite Innovation Keynote," Microsoft's Chief Executive Officer, Mr. Satya Nadella, boasted of the numerous benefits Microsoft had achieved by incorporating FPGAs into its cloud computing offerings: "We now have FPGA support across every compute node of Azure. That means we have the ability, *through the magic of the fabric* that we have built, to distribute your machine learning tasks, your deep neural nets, to all of the silicon that is available, so that you can get that performance, that scale." Exhibit 31 at 9 (emphasis added). The "magic of the fabric" did indeed allow Microsoft to scale

¹ <https://www.statista.com/statistics/505248/worldwide-platform-as-a-service-revenue/>

up an FPGA-based cloud computing solution. But the “magic of the fabric” was not invented by Microsoft, it was copied from ThroughPuter.

8. According to Dr. Doug Burger (with whom ThroughPuter corresponded about its technology), the “magic of the fabric” referred to by Microsoft’s CEO, use of FPGAs “allows us to do things on a scale that hasn’t been done before.” According to Dr. Burger: “It gives us the most powerful cloud, the most flexible cloud, and the most intelligent cloud.” *Id.* at 30, 34.

9. The loss of ThroughPuter’s technical competitive advantage through Microsoft’s copying based infringement, which Microsoft admits gives it on the order of 100-fold performance and efficiency gain, has devastated ThroughPuter’s business, including the ability to raise sufficient startup capital, gain collaborators, partners and initial customers, and to generally enter the market with a differentiated or more cost-efficient solution. ThroughPuter has survived notwithstanding Microsoft’s anticompetitive tactics due to its ability to innovate in the fields of application of its core technology.

10. ThroughPuter now brings this action based on Microsoft’s willful infringement of the patents that form the causes of action herein.

NATURE OF THE ACTION

11. This is an action for patent infringement arising under the Patent Laws of the United States, 35 U.S.C. § 1 *et seq.*, including 35 U.S.C. § 271.²

12. ThroughPuter brings this action to halt Microsoft’s infringement of its rights under the Patent Laws of the United States, 35 U.S.C. § 1 *et seq.*, which arise under the following patents:

² In addition to the patents asserted in this action, ThroughPuter owns certain patent applications directed to the technology at issue in this action pending under the following U.S. Patent Application Nos. 16/798,310; 16/834,961; 16/812,158; 16/894,177; 17/195,174; and 17/212,903.

- U.S. Patent No. 10,963,306 (“the ’306 patent”), which is attached hereto as Exhibit 1,
- U.S. Patent No. 10,620,998 (“the ’998 patent”), which is attached hereto as Exhibit 2,
- U.S. Patent No. 10,437,644 (“the ’644 patent”), which is attached hereto as Exhibit 3,
- U.S. Patent No. 10,430,242 (“the ’242 patent”), which is attached hereto as Exhibit 4,
- U.S. Patent No. 10,318,353 (“the ’353 patent”), which is attached hereto as Exhibit 5,
- U.S. Patent No. 10,310,902 (“the ’902 patent”), which is attached hereto as Exhibit 6,
- U.S. Patent No. 10,133,599 (“the ’599 patent”), which is attached hereto as Exhibit 7,
- U.S. Patent No. 9,632,833 (“the ’833 patent”), which is attached hereto as Exhibit 8, and
- U.S. Patent No. 9,424,090 (“the ’090 patent”), which is attached hereto as Exhibit 9.

THE PARTIES

13. Plaintiff ThroughPuter, Inc. is a Delaware corporation having a principal place of business within this District at 249 Richmond Road, Williamsburg, VA 23185 and is licensed to sell goods and services in the Commonwealth of Virginia and this District under Virginia Entity ID number 11115424. Plaintiff has been developing its technology platform based services in this

District since at least May 2020. Plaintiff owns over 50 issued domestic and foreign patents and pending applications protecting its products, services and technologies including those offered for sale within this District. ThroughPuter's President, Mark Sandstrom, is the named inventor on each of such patents and applications.

14. Plaintiff is a member of The College of William & Mary's Launchpad business incubator, also known as the Miller Entrepreneurship Center. Plaintiff's physical office space has been the William & Mary Launchpad business incubator since May 2020, although physical usage of that space has been limited due to the COVID-19 pandemic.

15. From March 2020 to present, ThroughPuter has employed at any given time from two to four hardware and software developers and architects as employees and one architect as an independent contractor. But for the pandemic, all of the developer and architect employees would have been physically working in this District. Due to the pandemic, however, it has been necessary to accommodate remote work. From March 2020 to present, at least two of the developers and architects were physically working within this District. After the pandemic, Plaintiff intends all or most full-time technical team members to work at office space located within this District.

16. Plaintiff has developed and continues to develop within this District various products and services, including i) Estimator™, a machine learning Application Specific Processor ("ASP")-as-a-service offering of the ThroughPuter PaaS project, and ii) Grafword™, an artificial intelligence ("AI") powered, graphical authentication service that is a pilot application of the Estimator™ machine learning microservice.

17. Estimator™ provides a streaming machine learning ("ML") microservice, to support AI applications in unpredictably changing operating environments. Estimator™ allows its prediction models and logic parameters to be adjusted continuously while the microservice is in

operation, such that its predictions will stay tuned-in to the prevailing reality of its operating environment, as that may evolve over time or even change abruptly. An International Search Report recently conducted by the International Search Authority under the Patent Cooperation Treaty concluded that this technology is patentable. A beta version of the Estimator™ application programming interface is currently commercially available for 3rd party developer subscription at www.estimatorlab.com.

18. Grafword™ provides graphic based high-security password generation and authentication, such that the level of authentication challenge is adjusted according to a level of deviation of a given user's online session attributes from what is expected for the given username. Grafword™ thus provides both high security as well as, for the authentic users, convenience in online authentication. An International Search Report recently conducted by the International Search Authority under the Patent Cooperation Treaty also concluded that this technology is patentable. A beta version of Grafword™ is used for Estimator™ account creation and login: <https://estimatorlab.com/landing>.

19. Both Estimator™ and Grafword™ have the potential to change the space in which they are offered due to the advantages provided by ThroughPuter's claimed inventions such as increased throughput and latency.

20. Microsoft is a corporation organized and existing under the laws of the State of Washington with its principal place of business at One Microsoft Way, Redmond, WA 98052.

21. Microsoft may be served with process through its registered agent Corporation Service Company, 100 Shockoe Slip, 2nd Floor, Richmond, VA, 23219 - 4100.

22. On information and belief, Microsoft has been registered to do business in the Commonwealth of Virginia under Virginia Entity ID number F1157421 since about October 1993.

23. On information and belief, Microsoft has had regular and established places of business in this judicial district since at least 1987.

24. Microsoft currently maintains regular and established places of business throughout this District, including corporate sales offices, retail store locations and data centers.

25. Microsoft has corporate sales offices within this District in Tysons Corner, VA and Reston, VA.

26. Microsoft's Azure Data Center is located within this District at 101 Herbert Drive, Boydton, VA 23917.

27. On information and belief, Microsoft's data center in Boydton, VA spans over 1,000,000 square feet.

28. On information and belief, Microsoft has had at least one additional data center in this District since at least 2014.

29. On information and belief, Microsoft employs at least hundreds of people within this District.

JURISDICTION AND VENUE

30. This is an action for patent infringement which arises under the Patent Laws of the United States, 35 U.S.C. § 1 et seq.

31. This Court has subject matter jurisdiction at least under 28 U.S.C. §§ 1331 and 1338.

32. Venue is proper in this District under 28 U.S.C. § 1400(b) because Microsoft has committed acts of infringement and has a regular and established place of business in this District.

33. This Court has personal jurisdiction over Defendant Microsoft pursuant to due process and/or Virginia's Long Arm Statute because Microsoft has committed and continues to

commit acts of patent infringement, including acts giving rise to this action, within the Commonwealth of Virginia and this District, and because Microsoft recruits Virginia residents, directly or through an intermediary located in this state, for employment inside or outside this state.

34. The Court's exercise of jurisdiction over Microsoft would not offend traditional notions of fair play and substantial justice because Microsoft has established at least the required minimum contacts with the forum.

35. Microsoft maintains regular and established places of business throughout Virginia and in this District, including the aforementioned corporate sales offices, retail store locations, and the Azure Data Center.

36. Microsoft has substantial business contacts within this District and has purposefully availed itself of the privileges and benefits of the laws of the Commonwealth of Virginia.

BACKGROUND

37. This case involves ThroughPuter's patented cloud computing, computing acceleration and related technologies, which were developed starting in 2010.

38. As of that time, advancements in computing technologies had generally fallen into two categories. First, in the field conventionally referred to as high performance computing, the main objective has been maximizing the processing speed of a given computationally intensive program running on dedicated hardware. In this field, speed was traditionally achieved by assigning a combination of separate parallel processors to all work on the same program simultaneously. Second, in the field conventionally referred to as utility or cloud computing, the main objective has been to most efficiently share a given pool of computing hardware resources among a large number of client application programs.

39. Thus, in effect, one branch of computing innovation has been seeking to effectively use a large number of parallel processors to accelerate execution of a single application program by parallelizing its processing across a maximum possible number of processors. At the same time, another branch of computing innovation has been seeking to share a single pool of computing capacity among a large number of application programs to optimize utilization of processing capacity. The former efforts pursue maximizing processing speed of a single program. The latter efforts pursue maximizing utilization of processing capacity.

40. As of the time of ThroughPuter's pioneering patent filings starting in 2011, there had not been major synergies between the effort to increase processing speed of a single program on the one hand, and maximizing processing capacity utilization on the other. Indeed, pursuing one of these traditional objectives often happened at the expense of the other, placing the two objectives in tension with each other.

41. For instance, while dedicating an entire parallel processor based (super) computer to each individual application would increase processing speed of the individual programs, it would also cause severely sub-optimal computing resource utilization, as much of the capacity would be idle much of the time. On the other hand, while seeking to improve utilization of computing systems by sharing their processing capacity among a number of applications would lead to enhanced resource utilization, it also tended to slow down processing of individual programs. As such, the overall cost-efficiency of computing was not improving as much as improvements toward either of the two traditional objectives would imply: traditionally, increases in processing speed came at the expense of system utilization efficiency, while overall system utilization efficiency maximization came at the expense of individual application processing speed.

42. The foregoing tension was exacerbated by the fact that even mainstream application performance requirements were increasingly exceeding the processing throughput achievable from a single CPU core, *e.g.*, due to the practical limits being reached on the CPU clock rates. This created an emerging requirement for intra-application parallel processing (at ever finer grades) even for mainstream programs in order to pursue satisfactory processing speeds, while these programs were to be increasingly hosted on cloud platforms where the processing resources would be shared among programs of multiple clients.

43. These internally parallelized and/or pipelined (concurrent) enterprise and web applications would ultimately be largely deployed on dynamically shared cloud computing infrastructure by entities such as Microsoft using the technologies patented, pioneered and promoted by ThroughPuter.

44. Given the foregoing, there existed a need as of 2011 for supporting a large number of concurrent applications on dynamically shared parallel processing resource pools. This then-existing need for a new parallel computing architecture could be met by a system that enabled increasing the speed of executing application programs (including through execution of a given application in parallel across multiple processor cores and/or using hardware accelerators) while at the same time improving the utilization of the available computing resources.

45. To address these problems, ThroughPuter developed hardware implemented dynamic resource management functionality including a scheduler, placer, inter-task communications and input/output system for use with multicore processor arrays dynamically shared among multiple concurrent applications, preferably to be deployed on FPGA processors. To that end, in this technology approach, the manycore processor array involves a fabric of reconfigurable cores that can be on-demand programmed to supply the needed mix or match of

hardware accelerators. ThroughPuter's technology provided a cloud computing solution that enables accelerated processing speeds across multiple application programs while at the same time optimizing processing resource utilization.

46. An exemplary embodiment of ThroughPuter's Dynamic Parallel Execution (DPE) Environment™ (sometimes referred to as "DPEE™" or "DPEE™ technology") is depicted in, for example, ThroughPuter's U.S. Patent No. 9,424,090. Exhibit 9. Referring to the figures and tables of the representative '090 patent reproduced below, this preferred embodiment includes an **array of core slots 120**. The processors, or cores, are coordinated and coupled together with a **hardware-based manycore fabric managed by the controller**, which allows the client program tasks to be dynamically placed on processor cores of many-core arrays while communicating with each other directly and securely.

Figure 1

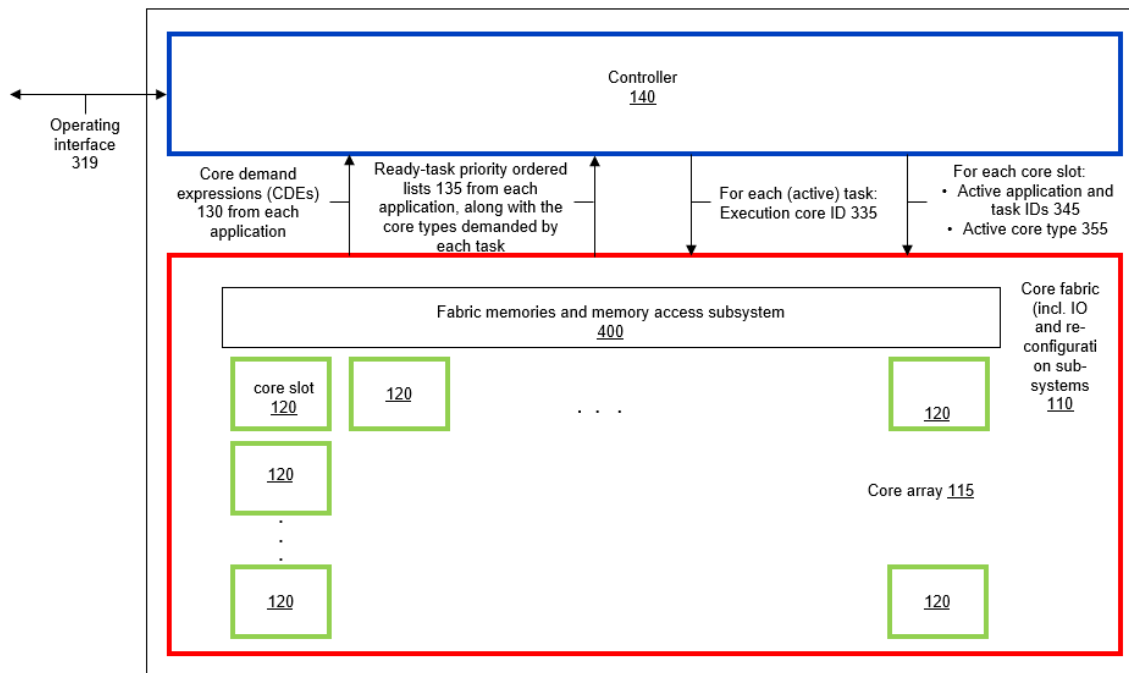


Exhibit 9, FIG. 1 (annotated).

47. ThroughPuter's hardware-based manycore fabric enables processing to be dynamically parallelized and hardware-accelerated, which results in optimized on-time processing throughput across the programs sharing an array of manycore processors. The effect for the client or end user of an accelerated service is increased processing speed and reduced cost base for delivering the application service, such that it becomes economically feasible for cloud service providers to support a range of performance intensive applications even without charge to end-users.

48. In a preferred embodiment, the **controller 140** monitors the processing load for each program sharing **the manycore array 110**, so that periodically re-optimized sets of tasks from the programs can be periodically assigned for execution on the **pool of processor cores 120**. Such periodic reassignment allows for optimal utilization of a pool of processing resources across a number of potentially competing, concurrent applications. ThroughPuter's exemplary embodiment enables assigning time-variable sets of application tasks for execution on the fabric of reconfigurable cores.

49. A variety of **application tasks 240** are assigned to be run on the **core fabric**, which includes **an array of core slots** within the **core fabric**. Each **application task** may be assigned to run on a **single processor core** within the **core fabric**. The cross-connect, which connects the various components of **the hardware-based manycore fabric**, is used by the **controller** to repeatedly optimize assignment of **tasks** to the appropriate **processor core slots**.

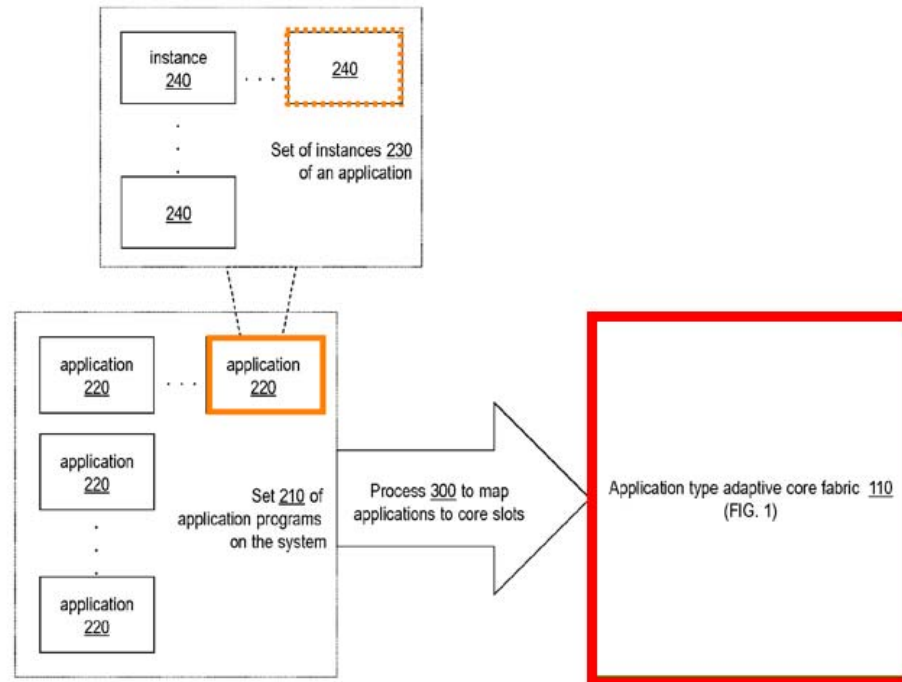
Figure 2

Exhibit 9, FIG. 2 (annotated).

50. The **controller** manages the **mapping** of **application tasks** to **processor cores** within **the core fabric**. The result is a system which dynamically and in real-time adapts to varying application processing loads to provide scalable, secure, high-performance and resource-efficient parallel cloud computing.

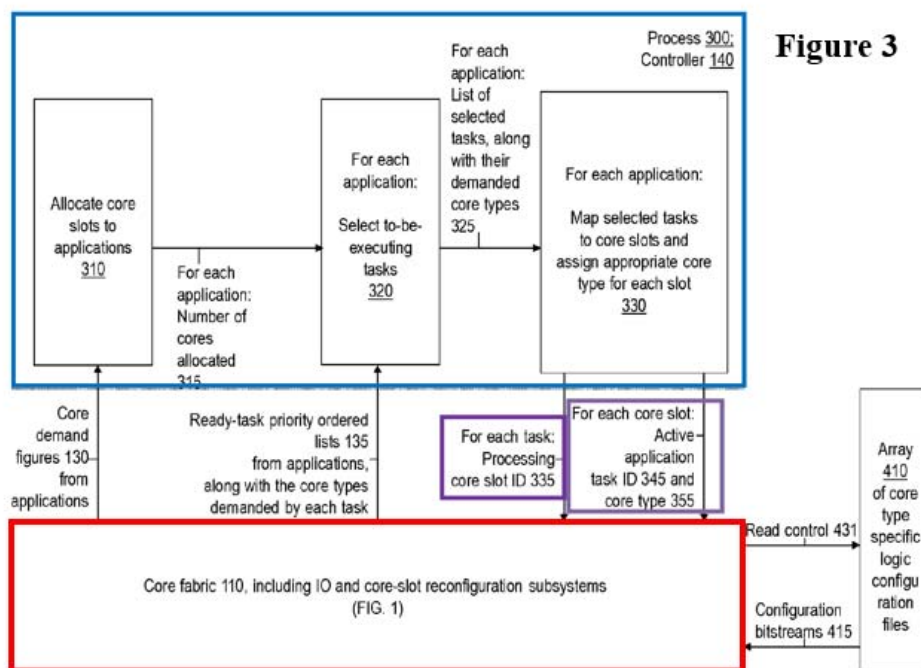


Exhibit 9, FIG. 3 (annotated).

TABLE 5

Core ID index	Application ID	Instance ID (within the application of column to the left)	Core type (e.g., 0 denotes CPU, 1 denotes DSP, 2 denotes GPU, 3 . . . 15 denotes an ASP for a particular function, etc.)
0	P	0	0
1	B	0	0
2	B	8	2
...
14	F	1	5
15	N	1	1

Exhibit 9, Table 5 (annotated).

51. This novel application load and type adaptive manycore fabric permits tasks or applications to be managed at a high degree of granularity with minimized processing overhead while providing various advantages not previously attainable.

52. In recognition of ThroughPuter's innovative achievements, ThroughPuter's Mark Sandstrom was invited to speak at various high performance and cloud computing conferences starting in 2012. GigaOm selected ThroughPuter as one of eleven finalists to present at Launchpad 2012 in San Francisco, CA. That same year, ThroughPuter was invited to present its Dynamic Parallel Execution Environment (DPEE) based PaaS approach at the high performance computing start-up showcase at the Supercomputing 2012 conference ("SC12") in Provo, UT.

53. In February 2013, Mr. Sandstrom reached out to Dr. Doug Burger, Director of Client and Cloud Applications at Microsoft Corporation via email indicating that ThroughPuter was looking for collaborators. The body of the email, which is attached as Exhibit 19, is reproduced below:

Doug,

I read at http://www.hpcwire.com/hpcwire/2013-02-01/kalray_produces_supercomputer-on-a-chip_for_industrial_applications.html your remarks at HiPEAC, concerning the need to modernize the computing architectures in view of the physical limits and future application requirements. In this context you may find ThroughPuter's cross-layer optimized platform architecture, based on dynamic parallel execution model, quite relevant -- please review <http://www.throughputer.com/platform.html>

ThroughPuter is looking for collaborators for the effort to make the advanced dynamic parallel program execution capabilities available for users (application developers) via PaaS model; possibly Microsoft Azure business unit might be interested in exploring the collaboration opportunities.

Feel welcome to share this call for collaboration among the appropriate parties at Microsoft, and naturally please get back to me for any questions etc. further discussions.

54. Dr. Burger responded "Thanks Mark ... I appreciate the note. I'll forward to the right people in Azure." *Id.*

55. In February 2013, the webpage at <http://www.throughputer.com/platform.html> had the following content at the time of the email exchange:

Platform Overview

ThroughPuter does not build on pre-cloud and sequential processing era concepts such as standalone processor cores and manycore processors as collections of them, or inter-core/process communications or operating systems retrofitted for parallel cloud computing. Instead, ThroughPuter [platform](#) architecture is cross-layer optimized for dynamic, high-performance, high-efficiency, secure cloud computing.

In ThroughPuter architecture, parallel processing hardware is not a mere manycore processor. And neither is its fabric network mere wires and switches between the cores. Nor is its operating system based on conventional sequential OS models with parallel processing as something of an afterthought.

ThroughPuter is a parallel program development and execution platform-as-a-service designed for dynamic cloud computing. For instance, ThroughPuter execution environment is an actual dynamic, secure cloud processor. In ThroughPuter execution environment, the hardware operating system, the adaptable fabric of cores, the fabric network and memory architecture as well as the contract management subsystems all are part of the integrated platform, and work seamlessly together. Collectively, the dynamic parallel processing platform achieves maximized processing throughput, per unit cost, across all client programs dynamically and securely sharing a pool of processing resources.

For illustration, please review technical overview of ThroughPuter PaaS solution:

<http://www.throughputer.com/uploads/ThroughPuterPaaSforParallelProcessing.pdf>

Quick Features -> Benefits -> Customer Value sheet:

<http://www.throughputer.com/uploads/DynamicParallelExecution-FeaturesBenefitsValue.pdf>

Hard-core technology paper on dynamic parallel execution environment of ThroughPuter PaaS:

<http://www.throughputer.com/uploads/DynamicParallelExecutionEnvironment.pdf>

56. The PDF entitled Parallel Program Execution in a Dynamically Shared Adaptive Manycore Processor is a document that closely follows the detailed description of ThroughPuter's '090 patent. Exhibit 20.

57. The PDF entitled Dynamic Parallel Execution via PaaS is a document that provides an overview of the innovations described in ThroughPuter's portfolio of patents and applications including the application leading to the '090 patent and the substantial advantages they provide. Exhibit 21.

58. The PDF entitled *Parallel Computing Development and Hosting Platform as a Service* is a presentation that provides an overview of ThroughPuter's DPPE. Exhibit 22. As shown in one of ThroughPuter's slides, reproduced below, the cores are dynamically re-assigned to instances or tasks of different applications, which are coded by color.

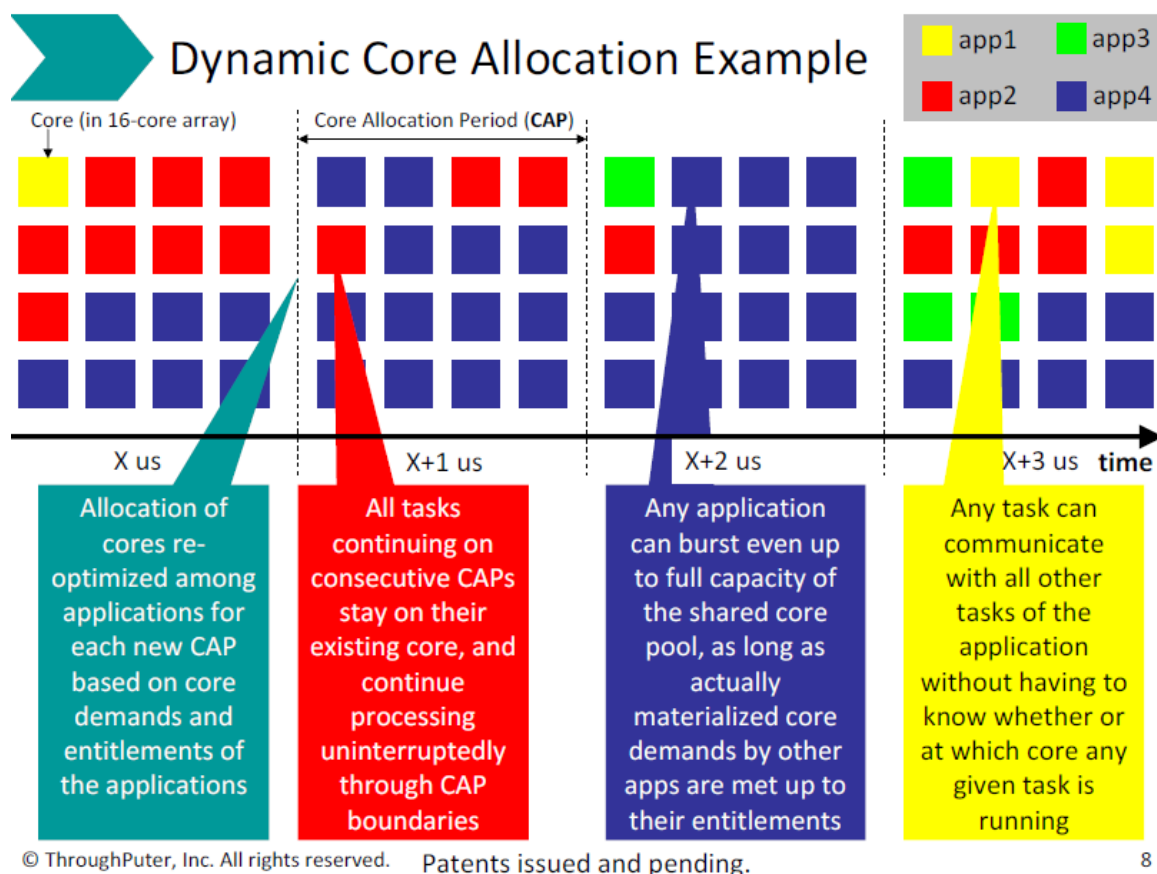


Exhibit 22 at 8.

59. Repeatedly, the pool of (sixteen) processing cores are each allocated among the four different application programs. In an illustrative example, these may be application programs used for **speech translation (blue)**, **encryption (red)**, **search ranking (yellow)** and **data compression (green)**. In the first cycle or period, the search ranking (**yellow**) application program is allocated one core while the encryption (**red**) program application and speech translation (**blue**) application program are allocated eight and seven cores, respectively. For the next period, several of the cores are reassigned to the speech translation (**blue**) program application to accommodate an increase in processing demand made by that application. In the third period, a further core is allocated to the speech translation (**blue**) program application and the data compression (**green**) program application is assigned one core, also in response to demand for those programs. In the

fourth period, the cores are more evenly allocated across all four applications because the spike in demand for speech translation has abated, or because of demand spikes of the other applications. In this way, the processing throughput and latency across all the applications sharing the given manycore processor array are optimized. The end result for the application programs is faster processing compared to what would be achievable for equal cost base under non-adaptive capacity allocation or without on-demand acceleration.

60. ThroughPuter's novel manycore fabric led to industry recognition of ThroughPuter and its technology. For example, in January 2013, ThroughPuter was invited to publish an article in the Cloud Computing Journal, discussing the PaaS based on novel manycore fabric. Exhibit 12.

61. In addition, in September 2014, Mr. Sandstrom presented at the FPGAworld conference in Stockholm, Sweden, on the topic of *Hardware Implemented Scheduler, Placer, Inter-Task Communications and IO System Functions for Manycore Processors Dynamically Shared among Multiple Applications*. Exhibit 23.

62. By the time of the 2014 FPGAworld conference, ThroughPuter had already been granted at least a dozen U.S. and United Kingdom patents protecting techniques enabling the advantages of its Dynamic Parallel Execution Environment™ (DPEE).

63. The following year, ThroughPuter was invited to present at the 2015 HPC Advisory Council Conference in Spain on the topic of executing multiple dynamically parallelized programs on dynamically shared cloud processors. A copy of the presentation is attached hereto as Exhibit 24.

MICROSOFT'S INFRINGING CLOUD COMPUTING ARCHITECTURE

64. Microsoft's infringing cloud computing platform is known as Azure which is the technology service platform accused of infringement herein.

65. In June 2015, Microsoft filed United States Patent Application Serial No. 14/752,800 (the “’800 application”). The corresponding patent is submitted herewith as Exhibit 11.

66. Microsoft’s Azure cloud platform includes or has included the functionality described in the ’800 application.

67. In its ’800 application, Microsoft explains that “graphs” or groups of FPGA processors (each designated with an H) are logically grouped together by a hardware layer fabric 106.

68. Each group of processors in Microsoft’s ’800 application can be dynamically reallocated and reconfigured to different application programs, such as **speech translation (blue)**, **encryption (red)**, **search ranking (yellow)** or **compression (green)**, as seen in Microsoft FIG. 15D below).

69. These reconfigurations are facilitated by the fact that the fabric provides addresses 1527 that identify the physical locations of each processing unit and the processing group (“graph”) to which a given processing unit is currently allocated.

70. The reconfigurations are performed, for instance, to accommodate increased demand for or from tasks of a given application program. In the system disclosed in the '800 application, service manager 1521 will repeatedly reallocate and reconfigure the processor groups or "graphs". In one example, as shown in Fig. 15D below, the service manager reallocates processor cores that were being used by the **encryption (red)** application (executed by graph 1503) to the **speech translation (blue)** application (executed by graph 1502). Compare Fig. 15C above with Fig. 15D below.

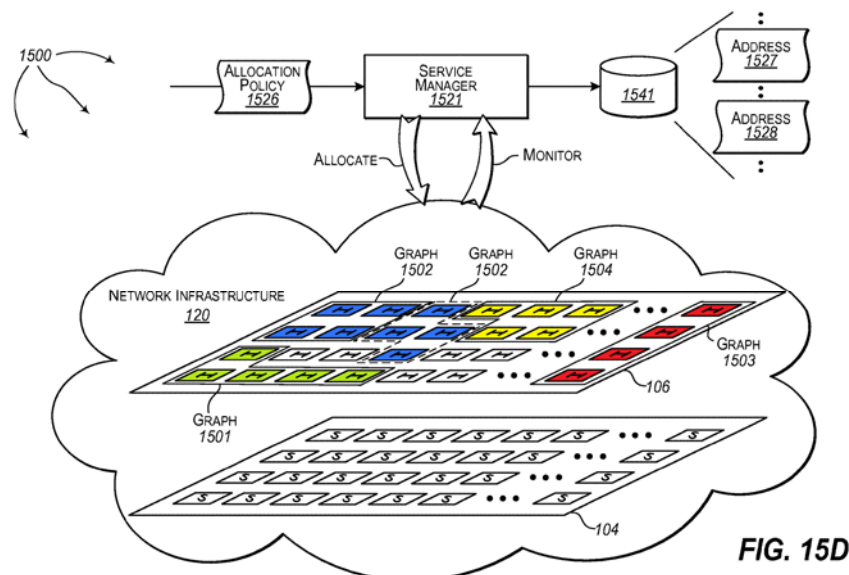
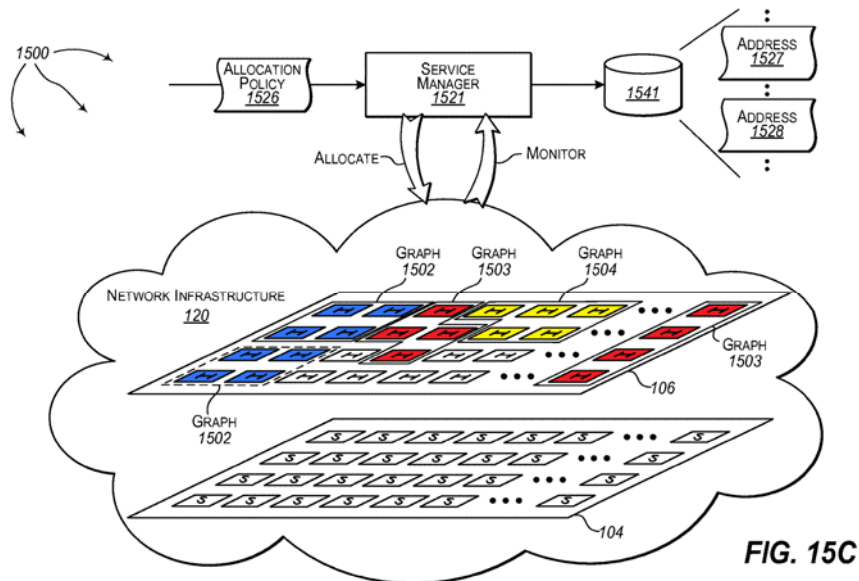


Exhibit 11, FIGs. 15C and 15D (annotated)

71. For this reallocation, the FPGA processor cores are reconfigured to have the functionality associated with the tasks being processed by the **speech translation (blue)** application or service (executed by graph 1502).

72. The service manager simultaneously reallocates FPGA cores from the speech translation group to accommodate increasing demand for a different processing task, for example **data compression (green, graph 1501)** as shown in the illustration above.

73. For this reallocation period, the processor groups handling, for example, **search ranking (yellow, graph 1504)** are not reallocated.

74. As illustrated above, the fabric addresses 1527 are updated to reflect the new processing core allocations.

75. As can be seen in claim 1 of the patent which issued from the '800 application (U.S. Patent No. 10,270,709, or the '709 patent, attached as Exhibit 11), Microsoft claims the functionality discussed, *supra*, at paragraphs 67-74.

76. On information and belief, Azure currently uses the technology described in its '709 patent and the Microsoft publications cited herein to manage at least a portion of its Azure workloads.

77. On information and belief, at least a portion of those workloads are handled at one or more of Microsoft's data center complexes within this District.

78. Claim 1 of the '709 patent is directed to a system that includes at least two sets of hardware acceleration components with each one capable of accelerating a different task or service as well as a processor and memory that maintains the mapping of each physical core address and

group assignment and reallocates the groupings based on variation in demand for a given group or core type.

79. Microsoft's **application tasks (S) such as search ranking** serves the same or similar purpose as the **application tasks 240** in ThroughPuter's representative '090 patent, discussed above. That is, both represent tasks to be performed on a hardware based processing system.

80. Microsoft's **hardware accelerators (H)** achieve the same or similar function as the **reconfigurable core slots 120** in ThroughPuter's representative '090 patent. That is, both represent reconfigurable processors on which the tasks can be performed.

81. **Microsoft's service managers 1521** (in coordination with the Resource Manager and the FPGA Managers) perform the same or similar function as the **controller 140 and/or the process 300**, in an embodiment of ThroughPuter's representative '090 patent, discussed above. That is, both assign tasks to an allocated group of reconfigurable processors.

82. Microsoft's **addresses 1527** serves the same or similar purpose as the **mappings shown in Figure 3 and Tables 4-5** of ThroughPuter's representative '090 patent, discussed above. That is, both provide addresses for the allocated processors to which the tasks are assigned.

83. The substantial identity between Microsoft's '709 patent and ThroughPuter's invention can be further appreciated from a side-by-side comparison of the patent claims granted to Microsoft and ThroughPuter. The column on the left shows claim 1 of the '306 patent, which claims priority to applications dating back to 2012. The column on the right shows the text of dependent claim 5 (including the limitations of independent claim 1) of Microsoft's '709 patent, which claims priority to an application filed in 2015. As can be appreciated from this side-by-side comparison, Microsoft obtained a patent on substantially the same technology taught by

ThroughPuter's patent applications. However, ThroughPuter's '306 patent is entitled to a priority date that is at least three years earlier than Microsoft's.

ThroughPuter's U.S. Patent No. 10,936,306, Independent Claim 1	Microsoft's U.S. Patent No. 10,270,709, Dependent Claim 5 (including the limitations of Independent Claim 1)
A method for task-switching on a multi-user parallel processing hardware architecture comprising a plurality of reconfigurable logic-based processing units, the method comprising:	A method in a data center comprising a plurality of servers interconnected via a network, wherein each of the plurality of servers comprises at least one acceleration component, wherein the at least one acceleration component comprises a plurality of hardware logic blocks interconnected using reconfigurable interconnects, the method comprising:
linking, through a first set of inter-task communication paths of the multi-user parallel processing hardware architecture, a first set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a first multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a first program, wherein each of the first set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of first tasks of the first program corresponding to a respective processing stage of the first multi-stage program instance;	linking a first set of acceleration components corresponding to a first set of the plurality of the servers into a first graph configured to accelerate a first service, wherein each of the first set of acceleration components is programmed to perform any of a first set of roles corresponding to the first service;
linking, through a second set of inter-task communication paths of the multi-user parallel processing hardware architecture, a second set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a second multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a second program, wherein each of the second set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of second tasks of the second program corresponding to a respective processing stage of the second multi-stage program instance;	linking a second set of acceleration components corresponding to a second set of the plurality of the servers into a second graph configured to accelerate a second service, different from the first service, wherein each of the second set of acceleration components is programmed to perform any of a second set of roles corresponding to the second service;

<p>maintaining, in a storage, a first location for a first reconfigurable logic-based processing unit of the first set of reconfigurable logic-based processing units executing the first multi-stage program instance such that a first one or more users and/or programs are enabled to communicate directly with the first multi-stage program instance;</p>	<p>maintaining in a storage a first address for the first graph such that the first service can request hardware acceleration from the first set of acceleration components;</p>
<p>maintaining, in the storage, a second location for a second reconfigurable logic-based processing unit of the second set of reconfigurable logic-based processing units executing the second multi-stage program instance such that a second one or more users and/or programs are enabled to communicate directly with the second multi-stage program instance; and</p>	<p>maintaining in the storage a second address for the second graph such that the second service can request hardware acceleration form the second set of acceleration components; and</p>
<p>in response to an increased demand for the second program, reallocating, by a controller comprising software and/or hardware logic configured to implement a load-adaptive allocation policy, at least one processing unit of the first set of reconfigurable logic-based processing units, the reallocating resulting in</p> <p>(1) switching the at least one of the first set of processing units from performing a task of the plurality of first tasks to performing one task of the plurality of second tasks, wherein switching comprises matching a first programming configuration of the at least one of the first set of processing units to a programming configuration demanded by the one task, and</p> <p>(2) adjusting, in storage, at least one of the first location or the second location to enable, through the respective location, direct communication to the other multi-stage program instance of the first multi-stage program instance or the second multi-stage program instance;</p>	<p>in response to an increased demand for hardware acceleration from the second graph, based on an allocation policy, a service manager comprising computer-executable instructions:</p> <p>(1) switching, using configuration data and without loading a new image file, at least one of the first set of acceleration components from performing any of the first set of roles to performing any of the second set of roles, and</p> <p>(2) adjusting at least one of the first address or the second address</p>

wherein the load-adaptive allocation policy is configured to facilitate minimizing reconfiguring the plurality of reconfigurable logic-based processing units.	wherein each of the first set of acceleration components and the second set of acceleration components is included in a hardware plane, and wherein the allocation policy is configured to facilitate minimizing switching of roles among acceleration components in the hardware plane.
--	--

84. On information and belief, Microsoft believed at the time of filing and still believes that the subject matter claimed in the '709 patent is patent eligible under 35 U.S.C. § 101.

85. On information and belief, Microsoft believed at the time of filing and still believes that this subject matter is novel under 35 U.S.C. § 102 in view of the references disclosed to and considered by the United States Patent & Trademark Office ("USPTO") during examination of their '800 application.

86. On information and belief, Microsoft believed at the time of filing and still believes that this subject matter is non-obvious under 35 U.S.C. § 103 in view of the references Microsoft disclosed to and considered by the USPTO during examination of their '800 application.

87. In 2016, Microsoft published a paper entitled *A Cloud-Scale Acceleration Architecture* authored by, among others, Adrian Caulfield (hereafter "Cloud-Scale Acceleration Architecture"). Exhibit 10.

88. Microsoft Azure includes or has included the functionality as described in Cloud-Scale Acceleration Architecture.

89. In Cloud-Scale Acceleration Architecture, Microsoft described a configurable core fabric similar to the configurable core fabric described in Microsoft's '800 application.

90. Fig. 13 of Cloud-Scale Acceleration Architecture illustrates a Resource Manager (RM) that, together with the Service (SM) and the FPGA Managers (FM), performs the same or similar function as the service manager 1521 described in Microsoft's '800 application.

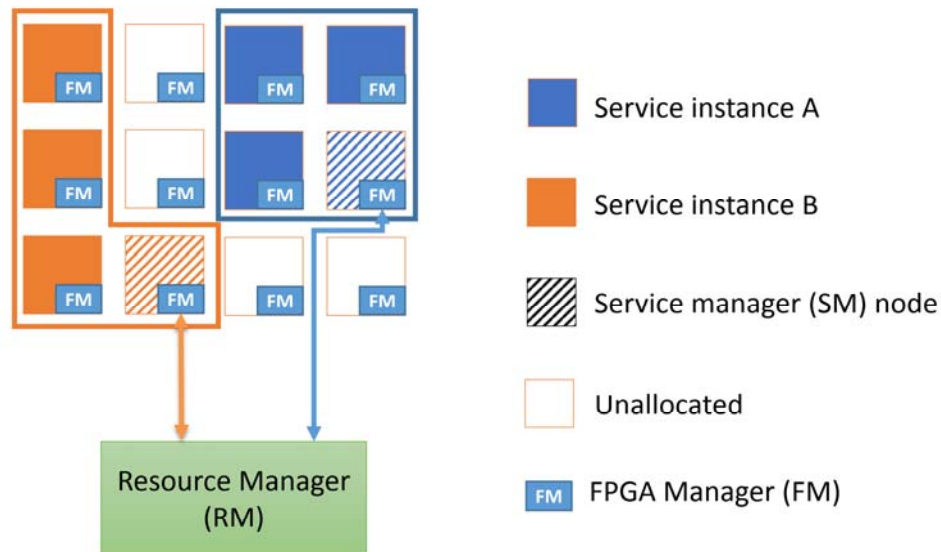


Fig. 13. Two Hardware-as-a-Service (HaaS) enabled hardware accelerators are shown running under HaaS. FPGAs are allocated to each service from Resource Manager's resource pool. Each service has a Service Manager node to administer the service on the allocated resources. Each FPGA has a lightweight FPGA Manager for managing calls from the Service Manager.

Exhibit 10, Fig. 13.

91. Each Service Manager (SM) described in Cloud-Scale Acceleration Architecture is associated with a given, dynamically sized and placed, FPGA processor group, each of which is referred to as a “service instance” in Cloud-Scale Acceleration Architecture.

92. Cloud-Scale Acceleration Architecture uses the phrase “service instance” to refer to the same or similar concept as the term “graph” in the '800 application.

93. Cloud-Scale Acceleration Architecture states that the functionality described therein is used to “accelerate . . . Azure infrastructure workloads,” among other use cases. Exhibit 10 at 13.

94. Microsoft's presentation entitled *Inside Microsoft's FPGA-Based Configurable Cloud* (Exhibit 15) demonstrates that Azure embodies or has embodied the functionality described

in Cloud-Scale Acceleration Architecture (Exhibit 10) and the Microsoft '709 patent (Exhibit 11).

95. As explained in the presentation, Microsoft Azure is a system comprised of a number of FPGA-accelerated servers, which process tasks from multiple application programs involving a controller that allocates FPGA processors from among a number of FPGA processors and then assigns application tasks to FPGA processors in order to pursue accelerated application processing while seeking efficient processing capacity usage.

96. A transcript and corresponding screenshots of Microsoft's presentation entitled *Inside Microsoft's FPGA-Based Configurable Cloud* are attached as Exhibit 15. The video is available at https://www.youtube.com/watch?v=v_4Ap1bjwgs.

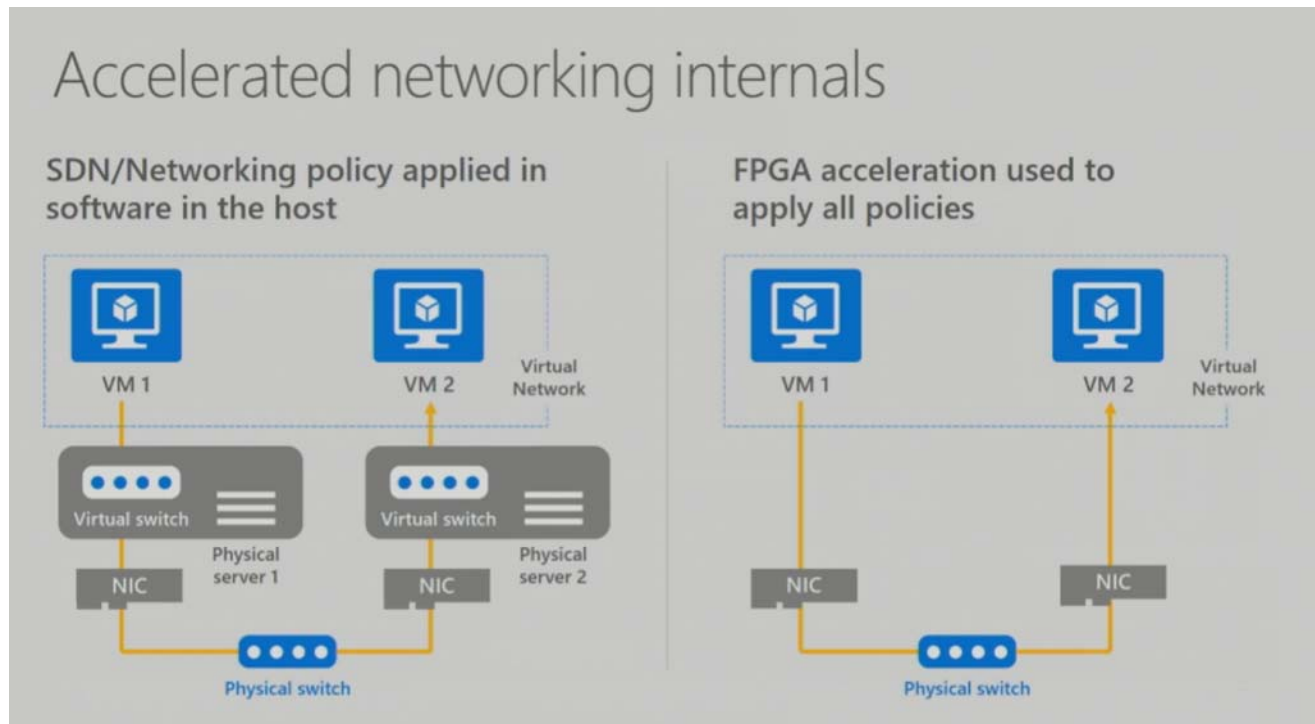


97. Microsoft Azure's RM interfaces with the Azure controller (AC) that applies policies to the data packets transmitted to and from the processors within Azure. Exhibit 15 at 39-42. The RM configures a hardware switch to execute the policies set forth in the hardware access control lists (ACLs), which list grant access rights to a given application program. *Id.*

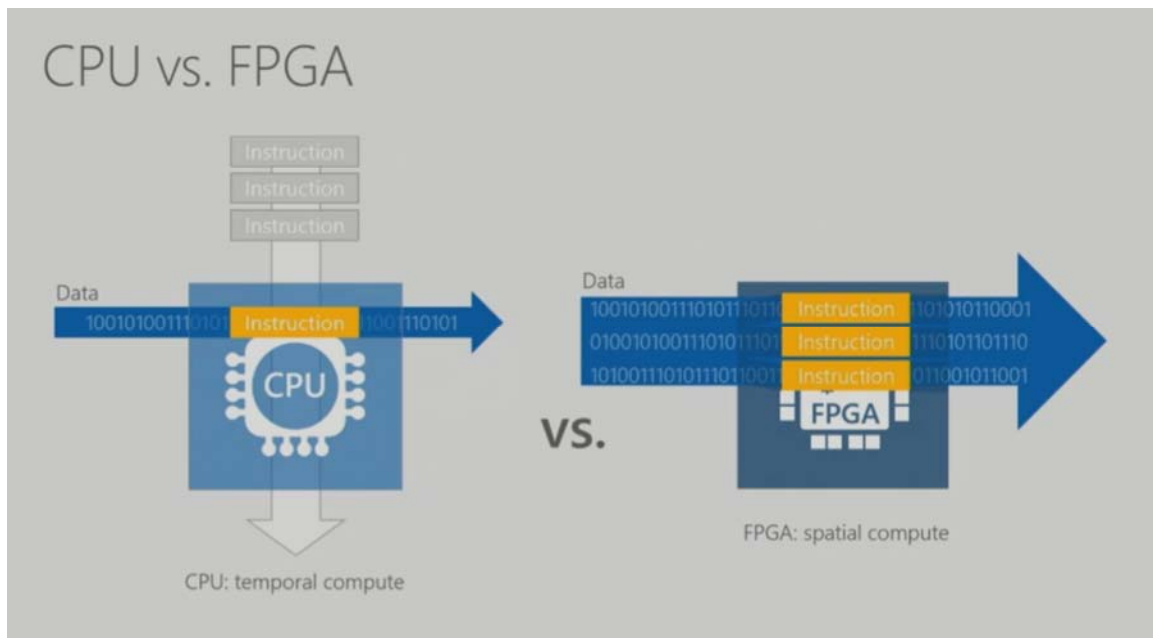
98. Under the direction of the AC, load balancing rules are practiced to pursue optimization of processing capacity by defining how tasks are distributed to FPGA Virtual Machines (*e.g.*, FPGA groupings or graphs in the parlance of Microsoft's '709 patent) (hereafter "FPGA Groupings"). Exhibit 15 at 43-50.

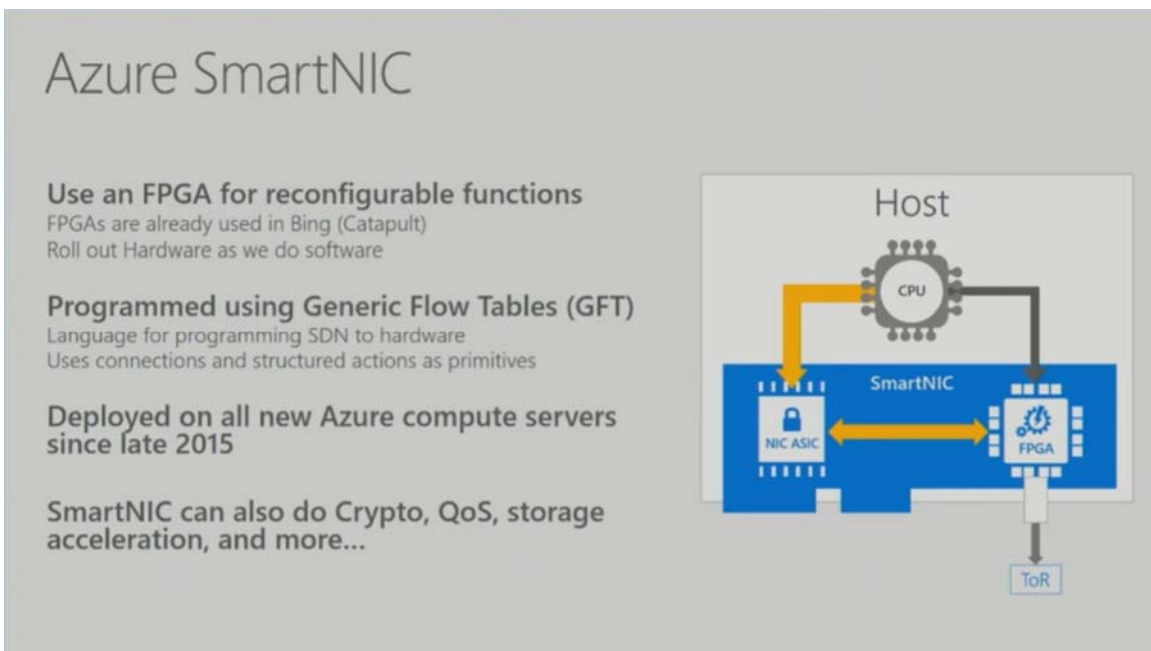
99. The AC instructs the fabric how to route each data packet associated with an application program or service instance task to a given processor based on the type of task to be performed and the available processing capacity on the Azure resource pool. The AC sets up load balancing rules in the load balancing control table (LB NAT) and pushes them out to all servers hosting FPGA Groupings.

100. Under the direction of the AC, Microsoft Azure's physical switch connects directly to the network interface cards (NICs) to route packets (*e.g.*, messages) to application or service instance tasks running on the FPGA Grouping. Exhibit 15 at 54-57. Each FPGA Grouping consists of several FPGA processors operating in a coordinated fashion under the control of the RM in conjunction with an SM, and a FPGA Manager or Node Manager (FM), as explained further below.



101. The FPGA processors are reconfigured to perform various functions for different application programs. *See, e.g.,* Exhibit 15 at 19-24, 54. As explained further below, the FPGA processors are dynamically reconfigured to be suited for a given task (*e.g.,* deep neural network (DNN) application tasks, SQL database tasks, and so on).

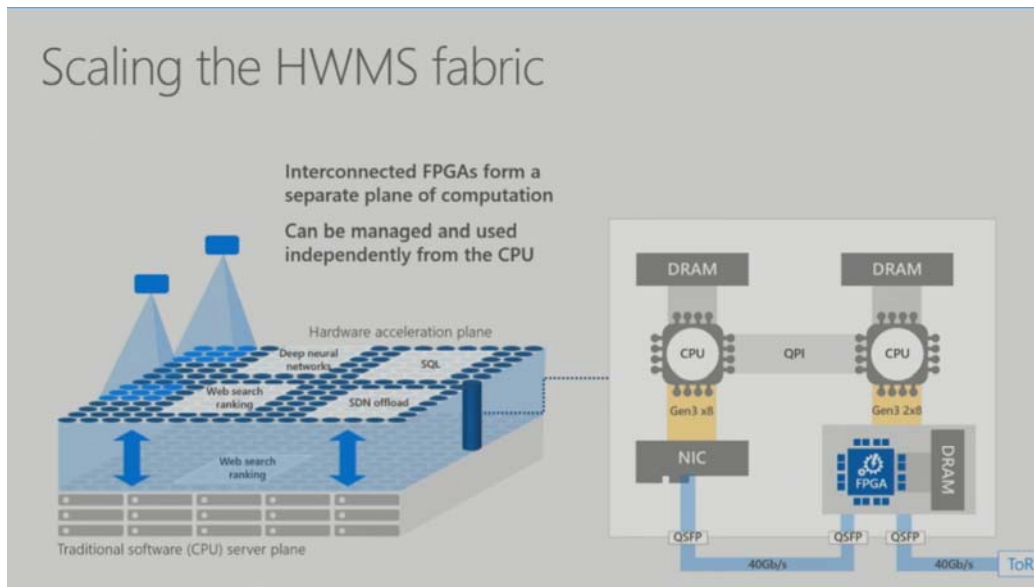




102. As shown in the above figure, an Azure Server includes an FPGA, a NIC, and a CPU. The FPGA is connected with the NIC and CPU, and the Azure Server is connected to other Azure Servers through a network, constituting a fabric of pooled FPGA accelerators. In an implementation disclosed in Fig. 8 in the '709 patent, an Azure Server includes an FPGA subdivided into separate configurable domains. In this implementation, the separate configurable domains are locally connected to each other through on-chip/on-board hardware. This forms a fabric of accelerator units that may be, at least in part, internal to a subdivided FPGA (hereafter the fabrics referred to in this paragraph are referred to as "Azure Fabric").

103. The FPGA processors are arranged in groups to execute tasks associated with various application programs. Exhibit 15 at 93-98. For example, one FPGA service instance or graph, *e.g.*, an FPGA Grouping, (the blue rectangle with downwardly projecting triangle) executes deep neural network (DNN) application tasks. *Id.* Another executes SQL database tasks, yet another executes network protocol offload tasks, and the last executes web search ranking tasks. Each FPGA Grouping includes several individual FPGA processors (illustrated as blue ovals in

the “hardware acceleration plane”).



104. In this manner, the AC repeatedly and dynamically rearranges task assignment to the array of processing units (*e.g.*, FPGA processors) while rearranging communication path connectivity for the array of processing units to optimize the application processing performance as well as the usage of processing capacity on the Azure system. Exhibit 15 at 95-99. The FPGA Groupings (and thus the associated connectivity) changes over time in reaction to the processing demand of the various application programs.

Benefits of Hardware Microservices

Decouple CPU to FPGA usage ratio

Flexibility: many services need a large number of FPGAs, others underutilize theirs

Deploy exactly as many instances as needed

Share accelerators (oversubscription)

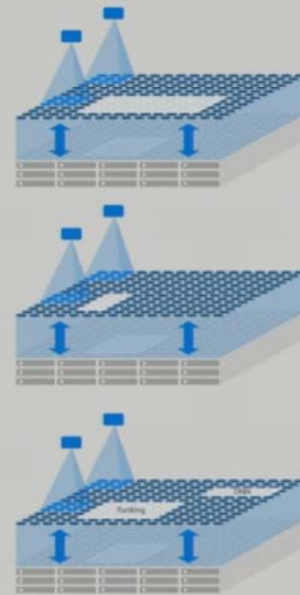
Many accelerators can handle load of multiple software clients

Consolidate underutilized FPGA accelerators into fewer shared instances

Increases efficiency & makes room for more accelerators

Expose multiple accelerators to a service

Many services need to access multiple types of accelerators



105. Microsoft is now using this FPGA fabric architecture in every new server it deploys in its data centers. Exhibit 28 at 8. Using this architecture, Microsoft obtained 40-100 fold performance improvements in Microsoft Bing's machine learning algorithms. *Id.* Microsoft has also announced the availability of Brainwave, an FPGA-based system for ultra-low latency deep learning for Azure. *Id.*

106. Using this architecture, Microsoft claims up to a 150-200 fold improvement in data processing throughput and up to a 50 fold improvement in energy efficiency. Exhibit 15 at 91, Exhibit 29. Latency has also lowered by about a factor of 10. Exhibit 15 at 86, Exhibit 29.

107. In Azure, the RM works in conjunction with both (a) the FM on each host processing unit, which keeps track of resource allocations on the FPGA processors, and (b) the SM, which assigns computing tasks to the FPGA processors. Exhibit 15 at 103-06.

108. In the implementation illustrated below, an SM in conjunction with the RM assigns one grouping of computing tasks of a ranking service application program (*e.g.*, Bing) to the FPGA

processors that have been allocated. Another grouping of FPGA processors is performing tasks for the Azure Data Link Analytics (ADLA) application, with those tasks having been assigned by another SM in conjunction with the RM, and so forth. In this way, the RM, SM and FMs ensure that each FPGA processor is receiving data from the appropriate queue and executing a task corresponding to the proper application. The RM, SM and FMs also ensure that the output of each FPGA accelerator is communicated to the proper output queue. The load balancing table is utilized to allocate application loads among FPGA processors and FPGAs among loads in order to pursue optimal usage of available processing capacity. In the illustrated example, the Ranking Service is expressing higher processing demand and/or has higher prioritization resulting in it getting allocated four FPGA cores and/or FPGA Groupings, compared with *e.g.*, ADLA that is assigned two FPGA accelerators.

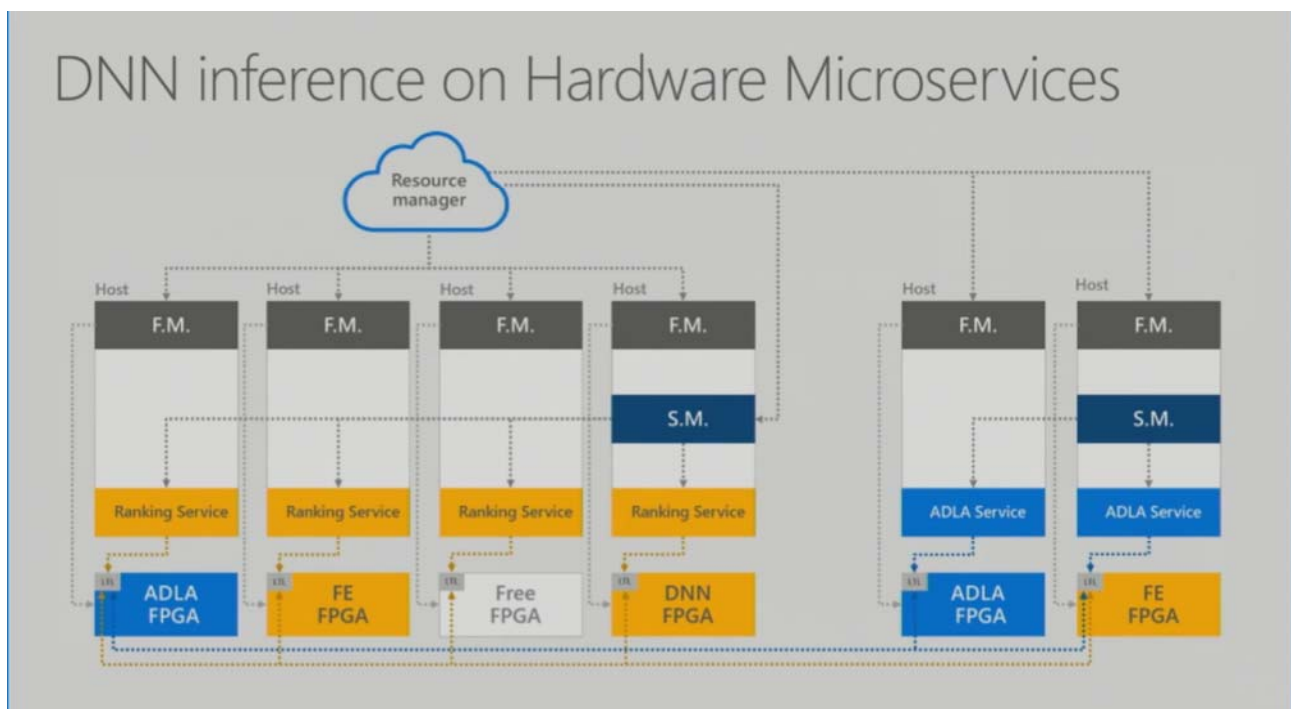


Exhibit 15 at 104.

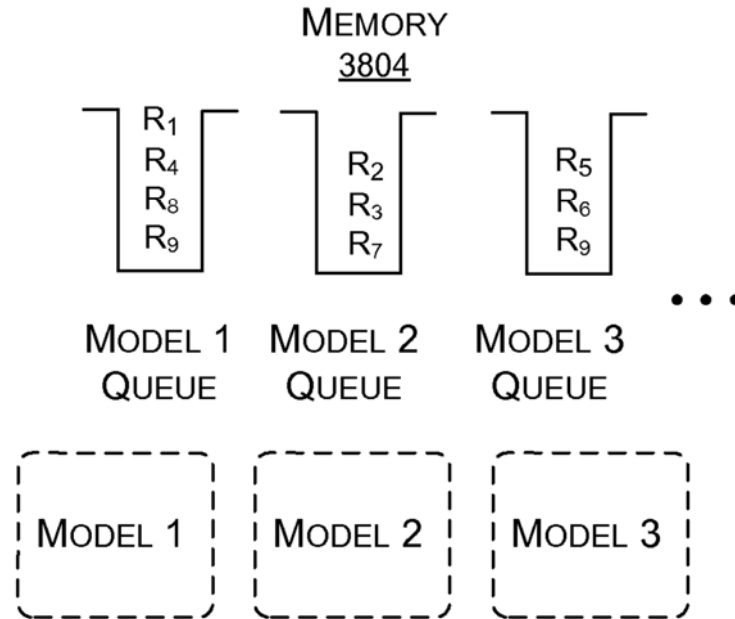
109. Azure's queuing of data as it moves between memories and the FPGA processors

is or has been accordant to descriptions in Microsoft's United States Patent No. 10,296,392 ("the '392 patent"). Exhibit 16.

110. On information and belief, Microsoft Azure embodies or has embodied the functionality substantially as described in the '392 patent.

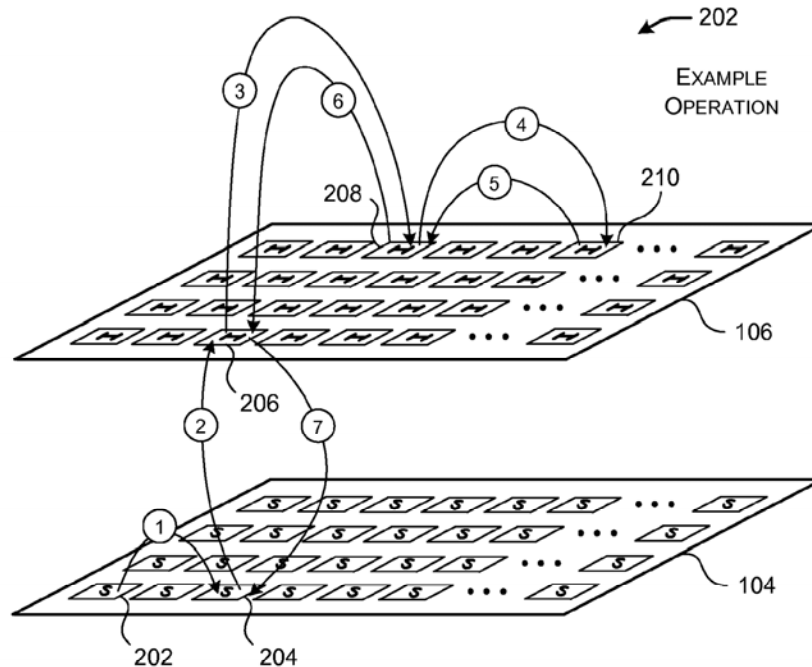
111. In connection with Fig. 38 (excerpted below), the '392 patent explains that different application programs such as a French language search engine, an English language search engine, and a German language search engine may have dedicated input buffers Model Queue 1, Model Queue 2 and Model Queue 3, respectively. Exhibit 16 at 34:62-35:34; see also Fig. 38. Data is fed from the queues to the FPGA processors based on policies such as queue fullness. *Id.* Generally, the queue manager (in coordination with at least the SM) will prioritize the processing of queries in the queue with the most queries. As an example, shown below, Model 1 Queue has queued four queries (R₁, R₄, R₈ and R₉) where Model 2 and Model 3 Queues have queued three queries. As such, Model 1 Queue would be prioritized above the other two.

112. On information and belief, Azure uses such queues to manage the transmission of data to, from and/or within the FPGA processors.



113. The distribution of tasks among the FPGA processors in Azure by the SM, RM, and FMs is also substantially described in the '392 patent.

114. In operations 1-6 illustrated in Fig. 2 of the '392 patent, each hardware acceleration component 208/210 (which may be an FPGA or a subdivision thereof) performs a task on the data and forwards the result to the next destination under the management of the head component of the FPGA Grouping (also referred to as the SM). The '392 patent further explains in connection with Fig. 36 that each acceleration component performs its respective function, such as mathematically combining feature values or compressing the data reflecting the feature values computed thus far.

**FIG. 2**

115. The scheduling and placement of tasks by Azure through the RM, SMs, and FMs is done taking into consideration a given task's readiness for execution. Exhibit 16 at 35:20-25; Exhibit 17 at pp. 287-88.

116. On information and belief, Azure embodies or has embodied the technology described in Exhibit 17.

117. Microsoft Azure also supports assured resource allocations for a client's FPGA Grouping for an additional fee. *See* Exhibit 18 at p. 3. Azure allocates FPGA processors based in part on the number of FPGA Groupings subscribed to by the client for its applications.

118. The foregoing functionality is further described in the Microsoft Azure documentation submitted herewith as Exhibits 12-14.

**MICROSOFT'S KNOWLEDGE OF
THROUGHPUTER AND ITS PATENTS**

119. Starting in 2011, ThroughPuter has invested in development and patent protection of its intellectual property consistently, in order to protect the competitive advantage to which it is entitled under the law as the inventor of the technology at issue in this case. Recognizing the promise of its technology and the benefits it could provide to Microsoft, ThroughPuter made consistent overtures to Microsoft in an effort to collaborate.

120. For example, as discussed above, starting in 2013, Mr. Sandstrom corresponded with Microsoft's Director of Client and Cloud Applications (Dr. Burger: a named inventor on the '709 patent) and others on Microsoft's cloud computing team concerning a potential collaboration between ThroughPuter and Microsoft. Exhibit 19.

121. At least through that correspondence, and no later than February 2013, Microsoft's cloud computing team was apprised of the technical details of ThroughPuter's Dynamic Parallel Execution Environment™ (DPEE).

122. At least through that correspondence, and no later than February 2013, Mr. Sandstrom provided Microsoft's cloud computing team a copy of the DPEE system description that corresponds with the preferred embodiments disclosed in ThroughPuter's '090 patent. *See* Exhibit 20.

123. From 2013 through 2016, ThroughPuter's President Mr. Sandstrom continued to correspond with Microsoft's cloud computing team concerning a potential collaboration between Microsoft and ThroughPuter.

124. On information and belief, Microsoft was made fully aware of the patent applications that issued as the '090 and '833 patents, starting at least as early as 2015.

125. In correspondence sent in 2015—to which Microsoft responded—ThroughPuter informed Microsoft of ThroughPuter’s patent applications leading to the ’090 and ’833 patents, as well as the Provisional Patent Application No. 61/934,747 providing comprehensive technical description for a reference implementation of the DPEE based manycore processor array technology. *See, e.g.*, Exhibit 32.

126. For example, in correspondence dated May 12, 2015, in response to a communication from ThroughPuter advising Microsoft of the patent applications that matured into the ’090 and ’833 patents, Dr. Burger stated that he had “forwarded [ThroughPuter’s] note onto the relevant people looking at any IP and they will contact you directly if interested. Doug.” Exhibit 30.

127. Shortly thereafter Microsoft filed its ’800 application (on June 26, 2015) naming Dr. Burger, to whom Mr. Sandstrom disclosed ThroughPuter’s technology, as an inventor. Following the filing of that application, the leader of Microsoft’s cloud computing team declined ThroughPuter’s offers for collaboration, saying in an August 3, 2015 email to Mr. Sandstrom: “I don’t think we can proceed ... we already have our plates very full. Thanks and best wishes.” Exhibit 32.

128. Thereafter, Mr. Sandstrom reminded Microsoft several times that the solution he shared was the subject of several U.S. and foreign patents.

129. In 2015, Mr. Sandstrom had a discussion with Microsoft representative, Dr. Derek Chiou, during which conversation, Dr. Chiou made clear that despite Microsoft’s knowledge of ThroughPuter’s technology and its patent protection, Microsoft would not be willing to explore any business relationship with ThroughPuter, including with respect to ThroughPuter’s intellectual property rights.

130. On May 16, 2016, Mr. Sandstrom informed Arun Justus at Microsoft that “ThroughPuter owns the IP rights for the key techniques that will be necessary in realizing any scalable solution for this must-solve challenge.” Exhibit 33.

131. In addition, on May 19, 2016, Mr. Sandstrom wrote Aki Siponen at Microsoft: “[All] cloud service providers will be facing the fundamental scalability solution . . . and ThroughPuter holds the key patented techniques that will be needed in delivering an effective solution to this must-solve challenge.” *Id.*

132. On May 17, 2017, via a LinkedIn message, Mr. Sandstrom provided notice to Dr. Burger of the issuance of the ’833 and ’090 patents, asserted herein under Counts 8 and 9. In response to Mr. Sandstrom’s message, Dr. Burger, stated “Thanks Mark, appreciate the nice message and will definitely review.”

133. In a 2017 Microsoft Research publication titled “The Feniks FPGA Operating System for Cloud Computing,” Microsoft describes its experimental development of FPGA processors with a hardware based operating system in highly similar terms as ThroughPuter in the above-referenced, earlier patent disclosures as well as publications such as the one presented by Mr. Sandstrom at FPGAworld in 2014. *Compare e.g., Exhibit 23 with Exhibit 25.*

134. In that 2017 publication, Dr. Chiou is credited with the “initial exploration which provides valuable experience for the system design.” (Exhibit 25 at p. 7, Acknowledgement). Dr. Chiou is the same individual with whom ThroughPuter spoke in 2015 after Microsoft had acknowledged ThroughPuter’s technology and its patent protection.

135. On January 2, 2018, Mr. Sandstrom sent a letter (the “January 2018 letter”) to Kevin Scott, Chief Technology Officer of Microsoft Corporation, informing Microsoft of the existence of the ’833 patent and other assets in ThroughPuter’s patent portfolio, and that ThroughPuter had,

at that time when it was lacking funds for commercialization of its technologies, “preference to sell the patents and the technology directly to an operating corporation such that enforcement of Throughputer’s intellectual property does not become necessary.”

136. On September 1, 2018, outside patent counsel for ThroughPuter sent Microsoft Chairman John W. Thompson a letter (“the September 2018 letter”) offering Microsoft the opportunity to license or acquire “ThroughPuter’s Dynamic Parallel Execution (DPE) technology, protected by a portfolio of forty-eight (48) patents issued and pending worldwide”. The September 2018 letter further explained that to the extent Microsoft’s current or planned products were built on certain cloud and enterprise computing technologies, “Microsoft needs to obtain a license from ThroughPuter for continued as well as planned future usage of [ThroughPuter’s] patents”.

137. The September 2018 letter specifically brought to Microsoft’s attention “claims 19 and 34 of U.S. Patent No. 9,632,833, claims 1 and 9 of U.S. Patent No. 9,424,090, and claims 1, 12 and 22 of U.S. Patent No. 10,061,615, along with a number of the patents’ dependent claims,” further noting “two Track 1 (expedited) patent applications, U.S. application serial nos. 16/014,658 and 16/014,674, each of which has received a notice of allowance.”

138. On September 1, 2018, ThroughPuter’s outside patent counsel sent the same letter referred to in the preceding paragraph to Microsoft’s then-General Counsel.

139. Microsoft did not respond to any of the letters referred to in the four preceding paragraphs.

140. Microsoft has not communicated with ThroughPuter since ThroughPuter advised Microsoft of the existence of U.S. Patent Nos. 9,424,090 and 9,632,833 in May 2017, which patents claim priority to patent applications filed in 2012.

141. On information and belief, Microsoft cut off communication with ThroughPuter in or around May 2017 because Microsoft realized that ThroughPuter had been issued patents with earlier priority dates on the technology approach underlying the Microsoft Azure PaaS and its development plans.

142. On information and belief, Dr. Burger and Microsoft's other personnel were instructed to stop communicating with ThroughPuter because Microsoft recognized the similarities between ThroughPuter's patent protected technology and the execution layer of the Azure PaaS and its development plans.

143. Rather than respect ThroughPuter's intellectual property rights, Microsoft made, upon information and belief, the deliberate decision to infringe ThroughPuter's patents. This conscious decision prevents ThroughPuter from being able to compete in the cloud computing space all while Microsoft used the enabling ThroughPuter technology to scale up the world's largest cloud-computing platform. ThroughPuter's patented technology allows Microsoft to claim up to 200-fold technical performance gains.

144. Microsoft recognized the novelty of ThroughPuter's technology by attempting to claim ThroughPuter's technology as its own. Specifically, as discussed herein, Microsoft sought and obtained patent protection on the fundamental technologies invented by ThroughPuter years earlier.

145. During the pendency of the patent application that matured into the '709 patent, Microsoft disclosed to the USPTO over two hundred fifty (250) documents that Microsoft represented as possible "prior art" to its alleged invention, but failed to disclose any information concerning ThroughPuter or its patents to the USPTO despite knowledge of the same.

146. During prosecution of Microsoft's '709 patent, Microsoft did not disclose to the USPTO ThroughPuter's '090 patent.

147. During prosecution of Microsoft's '709 patent, Microsoft did not disclose to the USPTO ThroughPuter's '833 patent.

148. During prosecution of Microsoft's '709 patent, Microsoft did not disclose to the USPTO any of the other ThroughPuter patent assets referenced in the materials provided to Microsoft by ThroughPuter.

149. During prosecution of Microsoft's '709 patent, Microsoft did not disclose to the USPTO any information concerning ThroughPuter or its DPEE technology.

150. To date, Microsoft has not disclosed to the USPTO in connection with any Microsoft patent application related to Azure any information concerning ThroughPuter or its DPEE, etc. technologies.

151. On information and belief, Microsoft was aware since at least 2017 that its Azure platform infringed ThroughPuter's patents, including the '090 and '833 patents.

152. Despite this knowledge, Microsoft has made the deliberate choice to refuse any business arrangement or resolution with ThroughPuter, opting instead to infringe ThroughPuter's patents, thereby preventing ThroughPuter from competing in the market ThroughPuter's technology has enabled.

153. On information and belief, Microsoft's infringement has been continuous, deliberate and in willful disregard of ThroughPuter's patent rights, including the '090 and '833 patents and the other patents identified in ThroughPuter's letters to Microsoft.

COUNT 1
(Infringement of U.S. Patent No. 10,963,306)

154. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

155. On March 30, 2021, the United States Patent and Trademark Office duly and legally issued the '306 patent, entitled "Managing Resource Sharing in a Multi-Core Data Processing Fabric." A copy of the '306 patent is attached as Exhibit 1.

156. Mark Sandstrom is the sole and true inventor of the '306 patent.

157. ThroughPuter, Inc. owns all right, title and interest to and in the '306 patent.

158. Microsoft infringes at least claims 1 and 4-7 of the '306 patent.

159. As demonstrated above in side-by-side comparison, claim 1 of ThroughPuter's '306 patent closely matches the claims of Microsoft's '709 patent. In other words, Microsoft later attempted to patent, and did obtain a patent on, the technology described by ThroughPuter's '306 patent disclosures, which were filed at least three years before the Microsoft '709 patent disclosure.

160. Claim 1 of the '306 patent is representative of the claims infringed by Microsoft and recites:

1. A method for task-switching on a multi-user parallel processing hardware architecture comprising a plurality of reconfigurable logic-based processing units, the method comprising:

linking, through a first set of inter-task communication paths of the multi-user parallel processing hardware architecture, a first set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a first multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a first program, wherein each of the first set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of first tasks of the first program corresponding to a respective processing stage of the first multi-stage program instance;

linking, through a second set of inter-task communication paths of the multi-user parallel processing hardware architecture, a second set of reconfigurable

logic-based processing units of the plurality of reconfigurable logic-based processing units into a second multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a second program, wherein each of the second set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of second tasks of the second program corresponding to a respective processing stage of the second multi-stage program instance;

maintaining, in a storage, a first location for a first reconfigurable logic-based processing unit of the first set of reconfigurable logic-based processing units executing the first multi-stage program instance such that a first one or more users and/or programs are enabled to communicate directly with the first multi-stage program instance;

maintaining, in the storage, a second location for a second reconfigurable logic-based processing unit of the second set of reconfigurable logic-based processing units executing the second multi-stage program instance such that a second one or more users and/or programs are enabled to communicate directly with the second multi-stage program instance; and

in response to an increased demand for the second program, reallocating, by a controller comprising software and/or hardware logic configured to implement a load-adaptive allocation policy, at least one processing unit of the first set of reconfigurable logic-based processing units, the reallocating resulting in

- (1) switching the at least one of the first set of processing units from performing a task of the plurality of first tasks to performing one task of the plurality of second tasks, wherein switching comprises matching a first programming configuration of the at least one of the first set of processing units to a programming configuration demanded by the one task, and
- (2) adjusting, in storage, at least one of the first location or the second location to enable, through the respective location, direct communication to the other multi-stage program instance of the first multi-stage program instance or the second multi-stage program instance;

wherein the load-adaptive allocation policy is configured to facilitate minimizing reconfiguring the plurality of reconfigurable logic-based processing units.

161. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '306 patent.

1. A method for task-switching on a multi-user parallel processing hardware architecture comprising a plurality of reconfigurable logic-based processing units, the method comprising:

162. The Microsoft Azure cloud platform (hereinafter, “Azure”) manages execution of a plurality of software application tasks on an array of reconfigurable logic-based processing units shared amongst multiple client programs. Azure is configured to, and Microsoft uses Azure to, switch application tasks across multiple processors within a group of processors.

163. At least the Azure Controller (AC) in conjunction with Resource Manager (RM), Service Managers (SM), and FPGA Managers (FM), is used to manage the assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs and/or FPGA processors.

linking, through a first set of inter-task communication paths of the multi-user parallel processing hardware architecture, a first set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a first multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a first program, wherein each of the first set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of first tasks of the first program corresponding to a respective processing stage of the first multi-stage program instance;

164. Azure links, through a first set of inter-task communication paths of the multi-user parallel processing architecture, a first set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a first multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a first program. In Azure, each of the first set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of first tasks of the first program corresponding to a respective processing stage of the first multi-stage program instance.

165. As described in Microsoft’s ’709 patent, the reconfigurable processing units

(acceleration components) are dynamically connected through hardware based communication paths that include a dedicated three port switch that works in coordination with the bypass control, multiplexers (connected to FIFO buffers 714, 716), the TOR interface, and TOR switch.

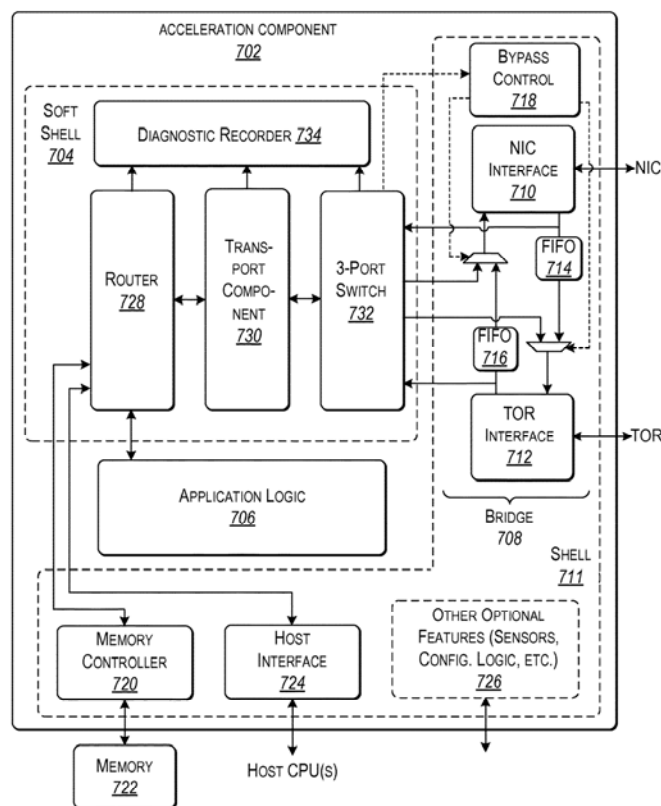


FIG. 7

Exhibit 11, FIG. 7.

166. Azure's SM, in coordination with the RM and FM, assigns groups of reconfigurable processing units such that some of the processing groups collectively execute a multi-stage program instance (*e.g.*, linking a set of processors to a multi-stage program instance) and individual processing units of an FPGA Grouping perform a single task associated with a stage of the multiple stage program instance. As illustrated by Fig. 5 of Microsoft's '709 patent, individual reconfigurable processing units (acceleration components H) perform their respective tasks corresponding to one of the stages and then communicates the results to the next processing unit

H, before the service or “graph” comprising said accelerators H outputs the result to the calling processing unit S.

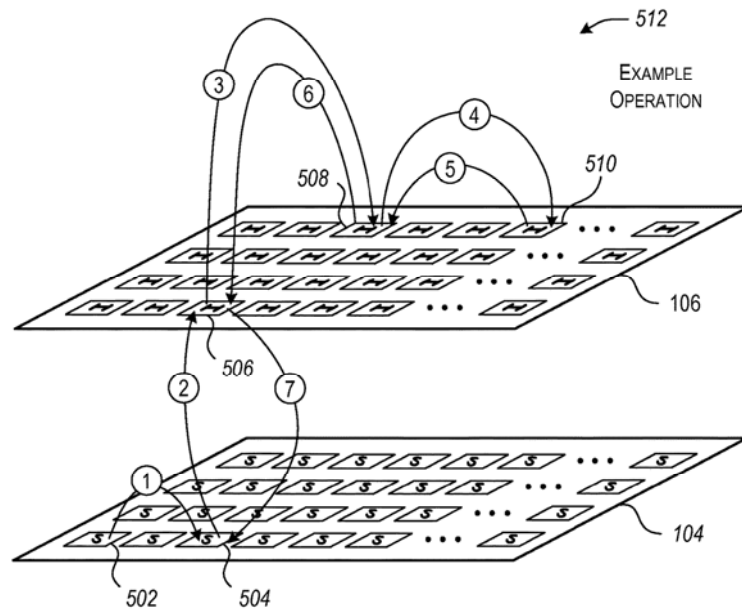


FIG. 5

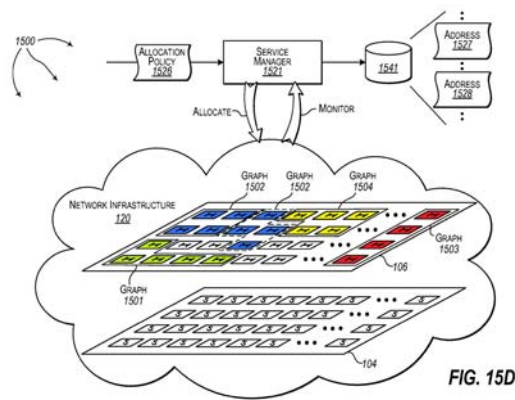
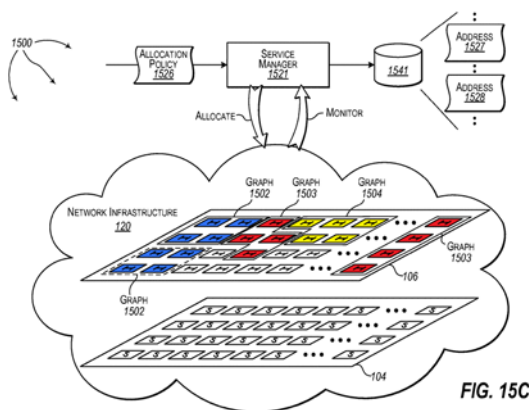
Exhibit 11, FIG. 5.

linking, through a second set of inter-task communication paths of the multi-user parallel processing hardware architecture, a second set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a second multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a second program, wherein each of the second set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of second tasks of the second program corresponding to a respective processing stage of the second multi-stage program instance;

167. Azure links through a second set of inter-task communication paths of the multi-user parallel processing hardware architecture, a second set of reconfigurable logic-based processing units of the plurality of reconfigurable logic-based processing units into a second multi-stage program instance configured to execute a corresponding set of interdependent processing stages of a second program. In Azure, each of the second set of reconfigurable logic-based processing units is programmed to perform a respective task of a plurality of second tasks of the

second program corresponding to a respective processing stage of the second multi-stage program instance.

168. Azure's AC works in conjunction with the RM, SMs, and FM to repeat the basic process described immediately above to assign a set of tasks of a second multi-stage program instance to a second group of reconfigurable processors. This is illustrated in Microsoft's '709 patent, which shows that different groups of reconfigurable processors H execute different tasks of a different multi-stage program instance during a given time period, shown in different colors in the annotated version of Fig. 15C below. During a subsequent time period, the arrangement of processors H has been reconfigured such that different sets of processors are performing the respective tasks, shown in different colors in the annotated version of Fig. 15D below.



169. The inter-task communication pathways between processing units vary between FPGA groups as well as over allocation time intervals, including because those groups (for instance graph 1503 and 1504) include different processing units H and because the same processing group (for instance graph 1502) includes varying processing units over time.

maintaining, in a storage, a first location for a first reconfigurable logic-based processing unit of the first set of reconfigurable logic-based processing units executing the first multi-stage program instance such that a first one or more users and/or programs are enabled to communicate directly with the first multi-stage program instance

170. Azure maintains, in a storage, a first location for a first reconfigurable logic-based processing unit of the first set of reconfigurable logic-based processing units executing the first multi-stage program instance such that a first one or more users and/or programs are enabled to communicate directly with the first multi-stage program instance.

171. Azure's RM and AC work in conjunction with at least the SMs to keep track of the addresses of each separately addressable FPGA processor group and maintain them in a storage that is accessible by the AC. According to Microsoft's '709 patent and as illustrated in Figs. 15C and 15D below, the service manager 1521 (which corresponds to the RM) stores addresses 1527-1528 for graphs 1502-1504 in storage 1541.

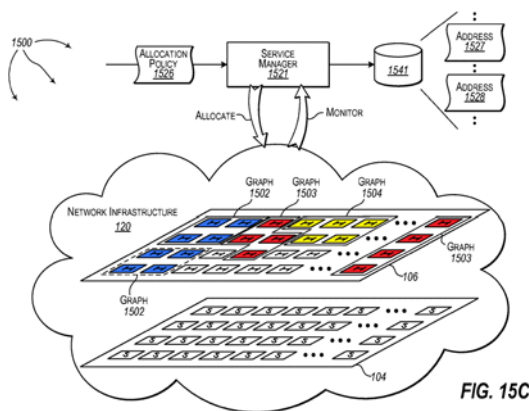


FIG. 15C

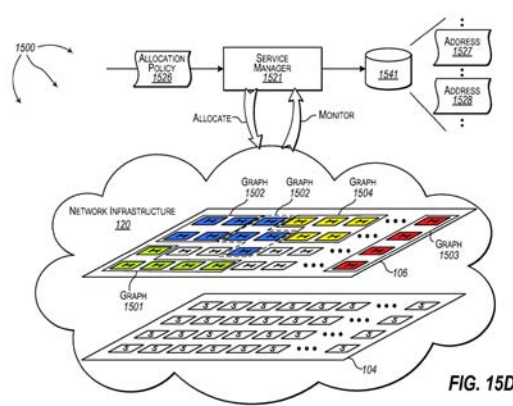


FIG. 15D

172. In Azure, client programs can communicate with the FPGA processor groups directly, which allow the client programs to communicate with the multi-stage program instance, including by contacting the RM and requesting that data be sent to a FPGA processor group and by receiving data output from a FPGA processor group. According to Microsoft's '709 patent, "[a]ddresses 1527, 1528, etc. stored in storage 1541 can be used to direct requests for acceleration

to the appropriate graphs.” Exhibit 11 at 19:57-59.

maintaining, in the storage, a second location for a second reconfigurable logic-based processing unit of the second set of reconfigurable logic-based processing units executing the second multi-stage program instance such that a second one or more users and/or programs are enabled to communicate directly with the second multi-stage program instance; and

173. Azure maintains in the storage, a second location for a second reconfigurable logic-based processing unit of the second set of reconfigurable logic-based processing units executing the second multi-stage program instance such that a second one or more users and/or programs are enabled to communicate directly with the second multi-stage program instance.

174. In Azure, the aforementioned processor group address storage and access functionality applies to and exists for each of the multiple FPGA processor groups or graphs.

in response to an increased demand for the second program, reallocating, by a controller comprising software and/or hardware logic configured to implement a load-adaptive allocation policy, at least one processing unit of the first set of reconfigurable logic-based processing units, the reallocating resulting in

175. In response to an increased demand for the second program, Azure reallocates by a controller comprising software and/or hardware logic configured to implement a load-adaptive allocation policy, at least one processing unit of the first set of reconfigurable logic-based processing units.

176. Azure’s RM and AC work in conjunction with at least the SMs to reallocate processors among the processor groups. According to Microsoft’s ’709 patent and as illustrated in Figs. 15C-D annotated below, the service manager 1521 coordinates the reallocation of the processors in graph 1503 to graph 1502.

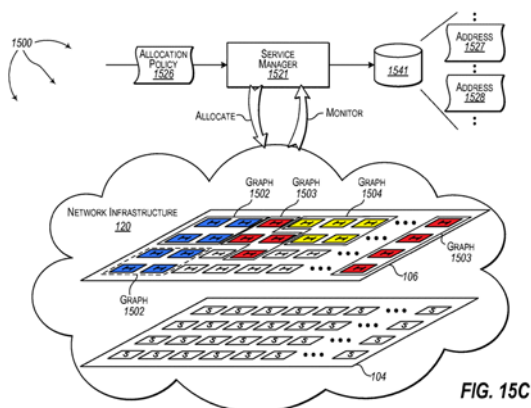


FIG. 15C

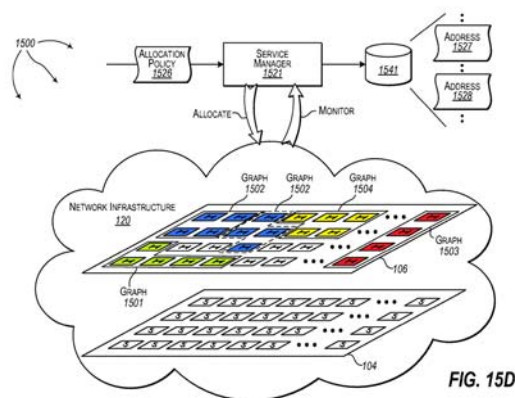


FIG. 15D

177. In Azure, this is done in response to a change in demand expressed by the application programs and/or their respective tasks, for instance by the accumulation of data at the input buffer for an application program and/or task. According to Microsoft's '709 patent, "[a]cceleration components [*i.e.*, reconfigurable processing units H] can be allocated in a manner that balances load in a hardware acceleration plane, minimizes role switching, and adapts to demand changes." Exhibit 11 at Abstract. The '709 patent further explains that "[d]uring continued monitoring, service manager 1521 can detect an increased demand for graph 1502. Based on the current allocation of acceleration components in hardware plane 106 and in view of allocation policy 1526, service manager 1521 can (re)allocate acceleration components to meet the increased demand for graph 1502." Exhibit 11 at 19:60-65.

(1) switching the at least one of the first set of processing units from performing a task of the plurality of first tasks to performing one task of the plurality of second tasks, wherein switching comprises matching a first programming configuration of the at least one of the first set of processing units to a programming configuration demanded by the one task, and,

178. Azure switches the at least one of the first set of processing units from performing a task of the plurality of first tasks to performing one task of the plurality of second tasks, wherein switching comprises matching a first programming configuration of the at least one of the first set

of processing units to a programming configuration demanded by the one task.

179. Azure's RM and AC work in conjunction with at least the SMs to reallocate and accelerate processors among the accelerator groups such that the need to reconfigure the accelerator logic blocks is minimized. For example, consider the situation in which a first calling application (*e.g.*, encryption) and a second calling application (*e.g.*, data compression) both preferentially request or require the same type of math function accelerator processor. In Azure, when the data compression application expresses increased demand for processing units relative to the encryption application, the RM and AC work in connection with the SMs to ensure that an accelerator logic block (from the encryption service) that is already configured to function as the math function accelerator is added to the data compression group instead of, *e.g.*, a logic block presently configured to function as a search ranking accelerator. This reduces the need to reconfigure the accelerator logic blocks, which would take additional time and impede overall efficiency.

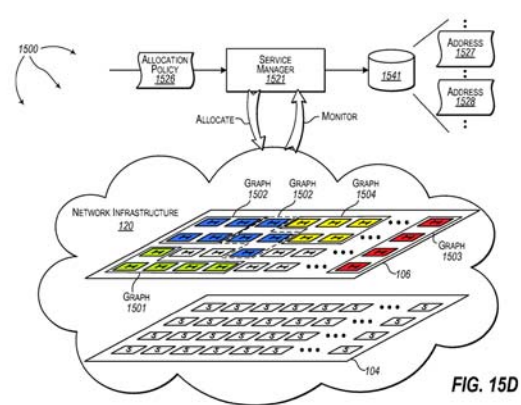
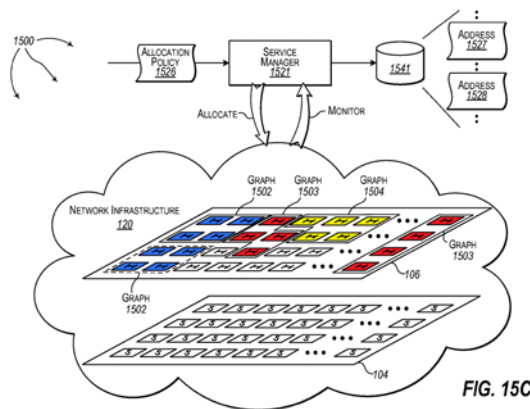
180. This functionality is described in Microsoft's '709 patent: "Configuration data can [be] used to switch between roles (or otherwise change behavior) at an acceleration component using the same underlying image [*i.e.*, processor configuration]. As such, roles (or other behavior) can be changed without having to load a new image file." Exhibit 11 at 20:26-30.

2) adjusting, in storage, at least one of the first location or the second location to enable, through the respective location, direct communication to the other multi-stage program instance of the first multi-stage program instance or the second multi-stage program instance,

181. Microsoft adjusts in storage, at least one of the first location or the second location to enable, through the respective location, direct communication to the other multi-stage program instance of the first multi-stage program instance or the second multi-stage program instance.

182. As discussed above, Azure's RM and AC work in conjunction with at least the SMs

to keep track of the addresses of each separately addressable processing group and record them in a storage that is accessible by the AC. As processing groups are modified in the manner depicted in Figs. 15C-D of Microsoft's '709 patent (annotated below), their corresponding addresses 1527-28 are updated in storage 1541.



183. According to Microsoft's '709 patent, “[t]he service manager maintains [*i.e.*, keeps updated] an address for the graph so that the service can request hardware acceleration from the group of interoperating acceleration components.” Exhibit 11 at 1:44-47.

184. As discussed above, and as illustrated by Fig. 5 of Microsoft's '709 patent, in Azure each reconfigurable processing unit (acceleration component H) performs a task corresponding to one of the stages and then directly communicates the results to the next processing unit H, before the service or “graph” comprising said accelerators H outputs the result to the calling processing unit S.

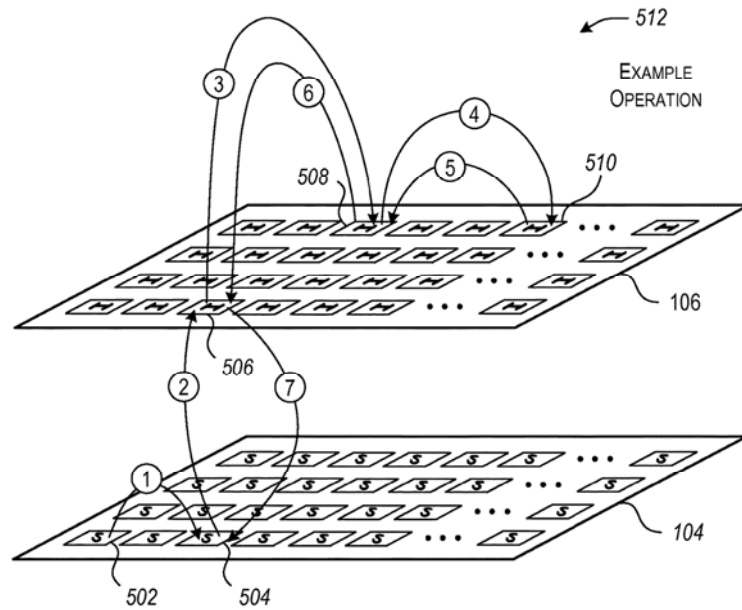


FIG. 5

wherein the load-adaptive allocation policy is configured to facilitate minimizing reconfiguring the plurality of reconfigurable logic-based processing units.

185. In Azure, the load-adaptive allocation policy is configured to facilitate minimizing reconfiguring the plurality of reconfigurable logic-based processing units.

186. As discussed above, Azure's RM and AC work in conjunction with at least the SMs to reallocate and reassign accelerator processors among the accelerator groups such that the need to reconfigure the accelerator logic blocks is minimized. For example, consider the situation in which a first calling application (*e.g.*, encryption) and a second calling application (*e.g.*, data compression) both preferentially request or require the same type of math function accelerator processor. In Azure, when the data compression application expresses increased demand for processing units relative to the encryption application, the RM and AC work in connection with the SMs to ensure that an accelerator logic block (from the encryption service) that is already configured to function as the demanded math accelerator is added to the data compression group

instead of, *e.g.*, a logic block presently configured to function as a search ranking accelerator. This reduces the need to reconfigure the accelerator logic blocks, which would take additional time and impede overall efficiency.

187. This functionality is described in Microsoft's '709 patent: "Configuration data can [be] used to switch between roles (or otherwise change behavior) at an acceleration component using the same underlying image [*i.e.*, processor configuration]. As such, roles (or other behavior) can be changed without having to load a new image file." Exhibit 11 at 20:26-30. The '709 patent further explains that "[a] service manager can minimize reconfiguration as much as possible when allocating a group of interoperating acceleration components to accelerate a service." Exhibit 11: at 7:14-17.

188. Upon information and belief, in accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the '306 patent no later than its issuance.

189. Microsoft continues, without license, to make use, offer for sale, import and/or sell in the United States services or products that infringe the '306 patent including specifically Microsoft Azure and its cloud computing functionalities.

190. Microsoft has directly infringed and continues to directly infringe the '306 patent by engaging in acts constituting patent infringement under 35 U.S.C. § 271(a) including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

191. Microsoft's continuing infringement of the '306 patent is knowing, intentional, and willful.

192. Microsoft has had knowledge of and notice of the chain of applications underlying

the '306 patent since at least January 2018 when ThroughPuter brought the '090 and '833 patents as well as the portfolio including a number of pending patent applications to the attention of Microsoft in the January 2018 letter, and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

193. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '306 patent under 35 U.S.C. § 271(b) and (c) through a range of activities.

194. First, on information and belief, Microsoft has induced infringement by, with knowledge of the '306 patent, controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '306 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

195. Second, on information and belief, Microsoft has, with knowledge of the '306 patent, induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '306 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

196. Third, on information and belief, Microsoft, with knowledge of the '306 patent,

has induced infringement by its customers through the creation and online posting of tutorial and “how-to” materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the ’306 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

197. Fourth, on information and belief, Microsoft has, with knowledge of the ’306 patent, induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers’ infringing use of the Azure platform.

198. Microsoft has engaged in the above activities with knowledge of the ’306 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

199. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the ’306 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the ’306 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft’s customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the ’306 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute at least a material part of the claimed invention of

the '306 patent.

200. Microsoft's infringement of the '306 patent has injured ThroughPuter in its business and property rights.

201. Microsoft's infringement of the '306 patent has been and is deliberate and willful and constitutes egregious misconduct. Upon information and belief, despite actual knowledge of the '306 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement.

202. Pursuant to 35 U.S.C. § 284, ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement in an amount to be determined at trial. Microsoft's infringement of the '306 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

203. Microsoft's infringement of the '306 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 2
(Infringement of U.S. Patent No. 10,620,998)

204. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

205. On April 14, 2020, the United States Patent and Trademark Office duly and legally issued U.S. Patent No. 10,620,998 ("the '998 patent"), entitled "Task Switching and Inter-Task Communications for Coordination of Applications Executing on a Multi-User Parallel Processing Architecture." A copy of the '998 patent is attached as Exhibit 2.

206. Mark Sandstrom is the sole and true inventor of the '998 patent.

207. ThroughPuter, Inc. owns all right, title and interest to and in the '998 patent.

208. Microsoft's infringes at least claims 1-7 and 9-21 of the '998 patent.

209. Claim 1 of the '998 patent is representative of the claims infringed by Microsoft

and recites:

1. A method for managing execution of a plurality of software applications on an array of processing units, the method comprising:

providing, for each software application of the plurality of software applications, one or more input buffers of a plurality of input buffers, each input buffer being provided for buffering processing load input directed to a respective software application of the plurality of software applications and being dedicated to the respective software application; and

repeatedly rearranging, by a controller comprising hardware logic and/or software logic, task assignment to the array of processing units and communication path connectivity for the array of processing units, wherein the array of processing units comprises a plurality of processing units of a reconfigurable type, each being configurable in two or more application specific configurations, and rearranging comprises, for each iteration of a plurality of iterations,

allocating, to each software application of at least a portion of a plurality of software applications as a set of active software applications, a number of units of the array of processing units at least in part in accordance with a plurality of demand expressions, each demand expression of the plurality of demand expressions corresponding to a different software application of the plurality of software applications, wherein each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of a processing load input at the one or more input buffers of the respective software application,

assigning, for each application of the set of active software applications, one or more task instances of the respective active application to one or more processing units of the plurality of processing units in accordance with the allocating,

wherein at least one task instance of the one or more task instances is a ready-to-execute task instance of the respective software application for processing the data at the one or more input data buffers of the respective software application,

each task instance of the one or more task instances is assigned to a different processing unit of the one or more processing units such that assigning results in each processing unit of the plurality of processing units being assigned only one respective task instance for any given iteration of the plurality of iterations, and

assigning comprises, for at least one task instance of the one or more task instances of a given software application of the plurality of software applications, reconfiguring the respective processing unit to an application specific configuration of the two or more application specific configurations, wherein the application specific configuration is associated with the at least one task instance, and

causing connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the plurality of processing units to connect respective processing load input to the respective software application that the respective processing load input is directed to.

210. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '998 patent.

1. A method for managing execution of a plurality of software applications on an array of processing units, the method comprising:

211. The Microsoft Azure cloud platform (hereinafter, "Azure") manages execution of a plurality of software applications on an array of processing units, *e.g.*, FPGA processors. Microsoft Azure is configured to, and Microsoft uses Azure to, manage and host a number of applications across multiple processors within a group of processors.

212. At least the AC in conjunction with RM, SMs FMs, is used to manage the assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs and/or FPGA processors.

providing, for each software application of the plurality of software applications, one or more input buffers of a plurality of input buffers, each input buffer being provided for buffering processing load input directed to a respective software application of the plurality of software applications and being dedicated to the respective software application; and

213. Azure provides, for each software application of the plurality of software applications, one or more input buffers of a plurality of input buffers in the form of queues.

214. A given input buffer is provided for buffering processing load input directed to a respective software application of the plurality of software applications and being specific to usage by the respective software application.

215. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, queues are associated with specific applications. Data from those queues are routed under the direction of the AC in conjunction with the RM, SMs and FMs to respective FPGA processors that have been at least temporarily or dynamically assigned to performing a task specific to an application program.

repeatedly rearranging, by a controller comprising hardware logic and/or software logic, task assignment to the array of processing units and communication path connectivity for the array of processing units, wherein the array of processing units comprises a plurality of processing units of a reconfigurable type, each being configurable in two or more application specific configurations, and rearranging comprises, for each iteration of a plurality of iterations

216. Microsoft Azure repeatedly rearranges, by a controller, which comprises hardware logic and/or software logic, task assignment to the array of processing units and communication path connectivity for the array of processing units, wherein the array of processing units comprises a plurality of processing units of a reconfigurable type, each being configurable in two or more application specific configurations, and rearranging comprises, for each iteration of a plurality of iterations.

217. Microsoft Azure includes a controller (*e.g.* at least the AC in conjunction with the RM, SMs, and FMs) comprising hardware logic and/or software logic that repeatedly rearranges task assignment to the array of processing units for each iteration.

218. The controller (*e.g.*, at least the AC in conjunction with the RM, SMs, and FMs)

also repeatedly rearranges communication path connectivity among the array of processing units.

219. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Microsoft’s Azure cloud platform includes an array of processing units that are of a reconfigurable type (*e.g.*, FPGA processors), each being configurable in two or more application specific configurations (*e.g.*, for DNN, SQL database programs, or a web search ranking algorithm or speech recognition, encryption or compression applications, demanding their respective forms of hardware acceleration).

allocating, to each software application of at least a portion of a plurality of software applications as a set of active software applications, a number of units of the array of processing units at least in part in accordance with a plurality of demand expressions, each demand expression of the plurality of demand expressions corresponding to a different software application of the plurality of software applications, wherein each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of a processing load input at the one or more input buffers of the respective software application,

220. Microsoft’s Azure (through *e.g.*, the RM in coordination with SMs, and FMs, which operate in conjunction with the AC) also allocates processing units to application programs at least in part in accordance with a plurality of demand expressions made by Azure that correspond to a different software application (*e.g.*, a program or group of programs that demand encryption acceleration).

221. In Azure, this is done in response to a change in demand expressed by the application programs and/or their respective tasks, for instance by the accumulation of data at the input buffer for an application program and/or task. According to Microsoft’s ’709 patent, “[a]cceleration components [*i.e.*, reconfigurable processing units H] can be allocated in a manner that balances load in a hardware acceleration plane, minimizes role switching, and adapts to demand changes.” Exhibit 11 at Abstract. The ’709 patent further explains that “[d]uring continued monitoring, service manager 1521 can detect an increased demand for graph 1502. Based on the

current allocation of acceleration components in hardware plane 106 and in view of allocation policy 1526, service manager 1521 can (re)allocate acceleration components to meet the increased demand for graph 1502.” Exhibit 11 at 19:60-65.

222. Each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of a processing load input at the one or more input buffers of the respective software application (*e.g.*, Azure allocates cores based on the measure of data to be processed for each software application such that the amount of processing capacity is allocated at least in part in accordance with the processing needs of the applications).

223. In Azure, each software application has its own processing needs, which are queued at the input buffers. The number of processor cores allocated to each application changes based on the applications’ need for processing capacity as reflected in the demand expressions. In the illustrative example depicted above, the Ranking Service program is allocated twice the number of cores that are assigned to the ADLA program due to the fact that the Ranking Service needs more processing capacity.

assigning, for each application of the set of active software applications, one or more task instances of the respective active application to one or more processing units of the plurality of processing units in accordance with the allocating, wherein at least one task instance of the one or more task instances is a ready-to-execute task instance of the respective software application for processing the data at the one or more input data buffers of the respective software application, each task instance of the one or more task instances is assigned to a different processing unit of the one or more processing units such that assigning results in each processing unit of the plurality of processing units being assigned only one respective task instance for any given iteration of the plurality of iterations, and

224. Microsoft Azure assigns (*e.g.* through the SM in coordination with the RM and FMs, which operate in conjunction with the AC), for each application of the set of active software applications, one or more task instances of the respective active application to one or more processing units of the plurality of processing units in accordance with the allocating of processors

from among the array of processing units.

225. At least one task instance of the one or more task instances, which are assigned to a processing unit, is a ready-to-execute task instance of the respective software application for processing the data at the one or more input data buffers of the respective software application.

226. In Azure, each task instance of the one or more task instances is assigned to a different processing unit of the one or more processing units such that assigning results in each processing unit of the plurality of processing units being assigned one respective task instance for any given iteration of the plurality of iterations.

assigning comprises, for at least one task instance of the one or more task instances of a given software application of the plurality of software applications, reconfiguring the respective processing unit to an application specific configuration of the two or more application specific configurations, wherein the application specific configuration is associated with the at least one task instance, and

227. Microsoft Azure's FPGA-enabled architecture performs assigning that includes, for at least one task instance of the one or more task instances of a given software application of the plurality of software applications, reconfiguring the respective processing unit to an application specific configuration of the two or more application specific configurations, wherein the application specific configuration is associated with the at least one task instance.

228. As explained above in the section entitled "Microsoft's Infringing Cloud Computing Architecture," the FPGA processors underlying the processor groups or graphs are dynamically reconfigured over time, in response to demand expressed by various applications, to perform specialized application-specific tasks such as web search ranking, data compression, DNN, or SQL database operations.

causing connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the

plurality of processing units to connect respective processing load input to the respective software application that the respective processing load input is directed to.

229. Microsoft Azure causes connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the plurality of processing units to connect respective processing load input to the respective software application that the respective processing load input is directed to.

230. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the queues such as message queues are associated with specific applications. Data from those queues are provided under the direction of the AC in conjunction with the RM, SMs, and/or FMs to respective FPGA processors that have been temporarily or dynamically assigned to perform a specific task for a specific application. Each FPGA processor performs a task such as data compression or web search ranking and returns the result back to the head unit or SM.

231. Upon information and belief, in accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the ’998 patent no later than its issuance.

232. Microsoft continues, without license, to make, use, offer for sale, import and/or sell in the United States services or products that infringe the ’998 patent including specifically Microsoft Azure and its cloud computing functionalities.

233. Microsoft has directly infringed and continues to directly infringe the ’998 patent by engaging in acts constituting patent infringement under 35 U.S.C. § 271(a) including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

234. Microsoft’s infringement of the ’998 patent has been, and continues to be knowing,

intentional, and willful.

235. Microsoft has had knowledge of and notice of the chain of applications underlying the '998 patent since at least January 2018 when ThroughPuter brought the '090 and '833 patents as well as the portfolio, including a number of pending patent applications, to the attention of Microsoft in the January 2018 letter, and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

236. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '998 patent under 35 U.S.C. § 271(b) and (c) through a range of activities.

237. First, on information and belief, Microsoft has induced infringement by, with knowledge of the '998 patent, controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '998 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

238. Second, on information and belief, Microsoft has, with knowledge of the '998 patent, induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '998 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the

Azure platform cloud computing services. Exhibit 26, Exhibit 27.

239. Third, on information and belief, Microsoft, with knowledge of the '998 patent, has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '998 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

240. Fourth, on information and belief, Microsoft has, with knowledge of the '998 patent, induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the Azure platform.

241. Microsoft has engaged in the above activities with knowledge of the '998 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

242. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '998 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the '998 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure

cloud platform, thereby infringing the '998 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute at least a material part of the claimed invention of the '998 patent.

243. Microsoft's infringement of the '998 patent has injured ThroughPuter in its business and property rights.

244. Microsoft's infringement of the '998 patent has been and is deliberate and willful and constitutes egregious misconduct. Upon information and belief, despite actual knowledge of the '998 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement.

245. Pursuant to 35 U.S.C. § 284, ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement in an amount to be determined at trial. Microsoft's infringement of the '998 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

246. Microsoft's infringement of the '998 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 3
(Infringement of U.S. Patent No. 10,437,644)

247. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

248. On October 8, 2019, the United States Patent and Trademark Office duly and legally issued the '644 patent, entitled "Task Switching and Inter-Task Communications for

Coordination of Applications Executing on a Multi-User Parallel Processing Architecture.” A copy of the ’644 patent is attached as Exhibit 3.

249. Mark Sandstrom is the sole and true inventor of the ’644 patent.

250. ThroughPuter owns all right, title, and interest to and in the ’644 patent.

251. Microsoft infringes at least claims 1-4, 6-7, 9-11, and 13-20 of the ’644 patent.

252. Claim 1 of the ’644 patent is representative of the claims infringed by Microsoft

and recites:

1. A system for managing execution of a plurality of software applications on an array of processing units, the system comprising:

a core fabric comprising

the array of processing units, and

a plurality of input buffers, each input buffer being provided for buffering processing load input directed to a respective software application of the plurality of software applications and being dedicated to the respective software application, wherein

each buffer of the plurality of input buffers is deployed in the core fabric apart from the array of processing units, and each software application of the plurality of software applications is provided one or more input buffers of the plurality of input buffers; and

a controller comprising hardware logic and/or software logic for performing operations for repeatedly reconfiguring task assignment to the array of processing units and communication path connectivity for the array of processing units, the operations comprising, for each iteration of a plurality of iterations,

allocating, to each software application of at least a portion of the plurality of software applications as a plurality of active software applications, a number of units of the array of processing units at least in part in accordance with a plurality of demand expressions, each demand expression of the plurality of demand expressions corresponding to a different software application of the plurality of software applications, wherein each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of the processing load input at the one or more input buffers of the respective software application, wherein allocating results in, for at least some of the plurality of software applications, a respective one or more allocated processing units of the array of processing units,

obtaining, for each software application of the plurality of active software applications, identification of one or more tasks of the respective software application,

assigning a respective task of the one or more tasks of each application of the plurality of active software applications to at least one respective processing unit of the respective one or more allocated processing units in accordance with the allocating, wherein

for at least a portion of the plurality of active software applications, the number of processing units allocated to the respective software application is less than a number of possible tasks of the respective software application available for assigning, and

assigning comprises, for each software application of at least a portion of the plurality of active software applications, identifying at least one activating task of the one or more tasks of the respective software application not assigned to the array of processing units for execution for a current iteration of the plurality of iterations, and

for each activating task of each software application of the portion of the plurality of active software applications, identifying at least one available unit of the array of processing units, each available unit corresponding to a respective deactivating task of at least one software application of the plurality of software applications, each deactivating task not assigned to the array of processing units for execution for a next iteration of the plurality of iterations, and assigning the at least one activating task to the at least one available unit, and

causing connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the respective one or more allocated processing units to connect respective processing load input to a task of the one or more tasks of the respective software application that the respective processing load input is directed to;

wherein, for consecutive iterations of the plurality of iterations and for at least one software application of the plurality of active software applications, a number of the one or more allocated processing units allocated to the respective software application varies with the measure of the processing load input at the one or more input buffers of the respective software application for the respective iteration, and the assigning comprises assigning a varying number of tasks of concurrent instances of the respective software application corresponding to variations in the number of processing units allocated to the respective software application between the consecutive iterations.

253. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '644 patent.

1. A system for managing execution of a plurality of software applications on an array of processing units

254. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Microsoft Azure comprises a system for managing execution of a plurality of software applications on an array of processing units. Microsoft Azure is configured to, and Microsoft uses Azure to, manage and host a number of applications across a pool of processing nodes. The AC in conjunction with the RM, FMs, and SMs is used to manage the assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs, and/or FPGA processors.

a core fabric comprising the array of processing units, and

255. Microsoft Azure includes a core fabric comprising an array of processing units. As discussed above, Azure includes the Azure Fabric, which is a parallel processing core fabric including an array of processing units.

a plurality of input buffers, each input buffer being provided for buffering processing load input directed to a respective software application of the plurality of software applications and being dedicated to the respective software application, wherein each buffer of the plurality of input buffers is deployed in the core fabric apart from the array of processing units, and each software application of the plurality of software applications is provided one or more input buffers of the plurality of input buffers; and

256. Microsoft Azure includes a plurality of input buffers, each input buffer being provided for buffering processing load input directed to a respective software application of the plurality of software applications and being dedicated to usage by a respective software application.

257. Each buffer of the plurality of input buffers is deployed in the core fabric apart from

the array of processing units, and each software application of the plurality of software applications is provided with one or more input buffers of the plurality of input buffers.

258. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, queues are associated with specific applications, and act as input buffers. Those input buffers are deployed in the Azure Fabric and are separate from the processors with which the buffers interface.

a controller comprising hardware logic and/or software logic for performing operations for repeatedly reconfiguring task assignment to the array of processing units and communication path connectivity for the array of processing units, the operations comprising, for each iteration of a plurality of iterations,

259. Microsoft Azure includes a controller comprising hardware logic and/or software logic for performing operations for repeatedly reconfiguring both task assignment to the array of processing units and communication path connectivity for the array of processing units.

260. Microsoft Azure includes a controller (*e.g.*, at least the AC in conjunction with the RM, SMs, and FMs) comprising hardware logic and/or software logic that repeatedly rearranges task assignment to the array of processing units.

261. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the controller also repeatedly rearranges communication path connectivity through the Azure Fabric among the array of processing units.

allocating, to each software application of at least a portion of the plurality of software applications as a plurality of active software applications, a number of units of the array of processing units at least in part in accordance with a plurality of demand expressions, each demand expression of the plurality of demand expressions corresponding to a different software application of the plurality of software applications, wherein each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of the processing load input at the one or more input buffers of the respective software application,

262. For each iteration, of a plurality of iterations, Microsoft Azure allocates, to each

software application of at least a portion of the plurality of software applications as a plurality of active software applications, a number of units of the array of processing units at least in part in accordance with a plurality of demand expressions, each demand expression of the plurality of demand expressions corresponding to a different software application of the plurality of software applications.

263. Microsoft's Azure cloud platform (through *e.g.*, the RM in coordination with the SMs and FMs, which work in conjunction with the AC) also allocates processing units to application programs at least in part in accordance with a plurality of demand expressions that correspond to the needs of a different software application (*e.g.*, a program for encryption services).

264. Each demand expression of the plurality of demand expressions is based at least in part on a measure of an amount of a processing load input at the one or more input buffers of the respective software application (*e.g.*, Azure allocates cores based on the measure of data to be processed for each software application such that the amount of processing capacity is allocated taking into consideration the respective processing needs of the software applications).

265. In Azure, each software application has its own processing needs, which are queued at the input buffers. The number of processor cores allocated to each application changes based on the applications' need for processing capacity as reflected in the demand expressions. In the illustrative example depicted above, the Ranking Service program is allocated twice the number of cores that are assigned to the ADLA program due to the fact that the Ranking Service needs more processing capacity.

wherein allocating results in, for at least some of the plurality of software applications, a respective one or more allocated processing units of the array of processing units, obtaining, for each software application of the plurality of active software applications, identification of one or more tasks of the respective software application

266. The allocating performed by Microsoft Azure results in, for at least some of the plurality of software applications, a respective one or more allocated processing units of the array of processing units, obtaining, for each software application of the plurality of active software applications, identification of one or more tasks of the respective software application.

267. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, a given application comprises a number of concurrent tasks to be executable by the processing units, which results in Azure identifying one or more tasks of the respective software applications.

assigning a respective task of the one or more tasks of each application of the plurality of active software applications to at least one respective processing unit of the respective one or more allocated processing units in accordance with the allocating, wherein

for at least a portion of the plurality of active software applications, the number of processing units allocated to the respective software application is less than a number of possible tasks of the respective software application available for assigning, and

268. Microsoft Azure assigns (e.g., at least the SMs in coordination with the RM and FMs, which work in conjunction with the AC) a respective task of the one or more tasks of each application of the plurality of active software applications to at least one respective processing unit of the respective one or more allocated processing units in accordance with the allocating, wherein for at least a portion of the plurality of active software applications, the number of processing units allocated to the respective software application is less than a number of possible tasks of the respective software application available for assigning.

269. As explained above in the section entitled “Microsoft’s Infringing Cloud

Computing Architecture,” at least the SMs in coordination with the RM and FMs, which work in conjunction with the AC, assign a number of processing cores based on the demand expressed by each application. In Azure, application programs are sometimes allocated fewer processing units than the number of queued tasks, including when other applications express higher demand.

assigning comprises, for each software application of at least a portion of the plurality of active software applications, identifying at least one activating task of the one or more tasks of the respective software application not assigned to the array of processing units for execution for a current iteration of the plurality of iterations, and for each activating task of each software application of the portion of the plurality of active software applications, identifying at least one available unit of the array of processing units, each available unit corresponding to a respective deactivating task of at least one software application of the plurality of software applications, each deactivating task not assigned to the array of processing units for execution for a next iteration of the plurality of iterations, and assigning the at least one activating task to the at least one available unit, and

270. Microsoft Azure (e.g., at least the SMs in coordination with the RM and FMs, which work in conjunction with the AC) identifies, for each software application of at least a portion of the plurality of active software applications, at least one activating task of the one or more tasks of a respective software application that was not assigned to the array of processing units for execution for a current iteration of the plurality of iterations.

271. For each activating task, Microsoft Azure identifies at least one available processing unit corresponding to a respective deactivating task of at least one software application of the plurality of software applications, with each such deactivating task being not assigned to the array of processing units for execution for a next iteration of the plurality of iterations.

272. Microsoft Azure assigns the at least one activating task to the at least one available unit.

273. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Azure changes the graphs or groupings of processors dynamically and

periodically (*e.g.*, by the RM in coordination with the SMs and FMs, which work in conjunction with the AC). In some situations, such as the scenarios described in the '709 patent, in a given reconfiguration some of the running tasks are deactivated from their group or application, such that a new task from potentially a different application can be activated and placed on an available processing unit.

causing connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the respective one or more allocated processing units to connect respective processing load input to a task of the one or more tasks of the respective software application that the respective processing load input is directed to;

274. Microsoft Azure causes connection, in accordance with the assigning, of the processing load input from each buffer of at least a portion of the plurality of input buffers to a different unit of the respective one or more allocated processing units to connect respective processing load input to a task among the one or more tasks of the respective software application that the respective processing load input is directed to.

275. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the queues such as message queues are associated with specific applications. Data from those queues are provided under the direction of the AC in conjunction with the RM, SMs, and FMs to respective processing units that have been temporarily or dynamically assigned to perform a task specific to an application program. When the data from one of the queues is connected to an FPGA processor of an FPGA Grouping, it may be a different processing unit than was previously connected to the queue.

wherein, for consecutive iterations of the plurality of iterations and for at least one software application of the plurality of active software applications, a number of the one or more allocated processing units allocated to the respective software application varies with the measure of the processing load input at the one or more input buffers of the respective software application for the respective iteration, and the assigning comprises assigning a varying number of tasks of concurrent instances of the respective software application corresponding to variations in the number of processing units allocated to the respective software application between the consecutive iterations.

276. In Azure, these steps are repeated and, for consecutive iterations of the plurality of iterations and for at least one software application of the plurality of active software applications, a number of the one or more allocated processing units allocated to the respective software application varies with the measure of the processing load input at the one or more input buffers of the respective software application for the respective iteration, and the assigning comprises assigning a varying number of tasks of concurrent instances of the respective software application corresponding to variations in the number of processing units allocated to the respective software application between the consecutive iterations.

277. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the RM in coordination with the SMs and FMs, which work in conjunction with the AC, allocates a number of processing cores based on the processing needs of each application. One task is assigned to each allocated processing core for a given allocation interval. During each periodic reallocation, the number of processing cores allocated to each application varies based on the processing needs of the various hosted applications.

278. On information and belief, in accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the ’644 patent no later than its issuance.

279. Microsoft has had knowledge of and notice of the chain of applications underlying the ’644 patent since at least January 2018 when ThroughPuter brought the ’090 and ’833 patents as well as its portfolio, including a number of pending patent applications, to the attention of

Microsoft in the January 2018 letter.

280. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '644 patent.

281. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '644 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

282. Microsoft's infringement of the '644 patent has been, and continues to be knowing, intentional, and willful. On information and belief, Microsoft has had knowledge of and notice of the application underlying the '644 patent since at least January 2018 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

283. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '644 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '644 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

284. Second, on information and belief, Microsoft has induced infringement by its

customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '644 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

285. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and “how-to” materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '644 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

286. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the Azure platform.

287. On information and belief, Microsoft has engaged in the above activities with knowledge of the '644 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

288. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '644 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the '644 patent, including at least the Azure platform as a whole and/or the individual

components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '644 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute a material part of the claimed invention of the '644 patent.

289. Microsoft's infringement of the '644 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '644 patent has been and is deliberate and willful and constitutes egregious misconduct. On information and belief, despite actual knowledge of the '644 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '644 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

290. Microsoft's infringement of the '644 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 4
(Infringement of U.S. Patent No. 10,430,242)

291. ThroughPuter repeats and realleges each and every allegation contained above as

though fully set forth herein.

292. On October 1, 2019, the United States Patent and Trademark Office duly and legally issued the '242 patent, entitled "Task Switching and Inter-Task Communications for Coordination of Applications Executing on a Multi-User Parallel Processing Architecture." A copy of the '242 patent is attached as Exhibit 4.

293. Mark Sandstrom is the sole and true inventor of the '242 patent.

294. ThroughPuter owns all right, title, and interest to and in the '242 patent.

295. Microsoft infringes at least claims 1-3, 6-7, 9-10, 12, 14-20, and 23 of the '242 patent.

296. Claim 1 of the '242 patent is representative of the claims infringed by Microsoft and recites:

1. A system for managing execution of a plurality of software applications on an array of processing units, the system comprising:
 - a core fabric comprising
 - the array of processing units, and
 - a plurality of input data buffers, each input data buffer being provided for buffering input data directed to a respective software application of the plurality of software applications and being dedicated to the respective software application, wherein
 - each buffer of the plurality of input data buffers is deployed in the core fabric apart from the array of processing units, and each software application of the plurality of software applications is provided one or more input data buffers of the plurality of input data buffers; and
 - a controller comprising hardware logic and/or software logic for performing operations for repeatedly reconfiguring task assignment to the array of processing units, the operations comprising, for each iteration of a plurality of iterations,
 - identifying, for each software application of at least a portion of the plurality of software applications, an amount of input data at one or more input data buffers of the plurality of input data buffers buffering data for the respective software application,
 - allocating, to each software application of the portion of the plurality of software applications, a number of processing units of the array of processing units based at least in part on the amount of input data buffered for the respective

software application, and for each software application of the portion,

- i) assigning one or more task instances of the respective software application for concurrent processing of the amount of input data to the number of processing units allocated to the respective software application by the allocating as one or more assigned instances, and
- ii) adjusting, based at least in part on a change in a count of units between the number of processing units allocated to the respective software application and a number of previously allocated processing units allocated to the respective software application during a previous iteration of the plurality of iterations, a relative portion of the amount of input data to be processed by at least one assigned instance of the one or more assigned instances;

wherein, for one or more iterations of the plurality of iterations where a current number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is increased by the allocating, adjusting comprises relatively decreasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application; and

wherein for one or more other iterations of the plurality of iterations where a present number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is decreased by the allocating, adjusting comprises relatively increasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application.

297. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '242 patent.

1. A system for managing execution of a plurality of software applications on an array of processing units, the system comprising:

298. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Microsoft Azure comprises a system for managing execution of a plurality of software applications on an array of processing units. Microsoft Azure is configured to, and Microsoft uses Azure to, manage and host a number of applications across a pool of processing nodes. The AC in conjunction with the RM, SMs, and FMs is used to manage the

assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs, and/or FPGA processors.

a core fabric comprising the array of processing units, and

299. Microsoft Azure includes the Azure Fabric, which is a core fabric comprising an array of processing units.

a plurality of input data buffers, each input data buffer being provided for buffering input data directed to a respective software application of the plurality of software applications and being dedicated to the respective software application, wherein each buffer of the plurality of input data buffers is deployed in the core fabric apart from the array of processing units, and each software application of the plurality of software applications is provided one or more input data buffers of the plurality of input data buffers; and

300. Microsoft Azure includes a plurality of input data buffers, each input data buffer being provided for buffering input data directed to a respective software application of the plurality of software applications and being dedicated by the respective software application, wherein each buffer of the plurality of input data buffers is deployed in the core fabric apart from the array of processing units, and each software application of the plurality of software applications is provided one or more input data buffers of the plurality of input data buffers.

301. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, queues are associated with specific applications, which act as input buffers. Those input buffers are deployed as part of the core fabric and separate from the processors with which the buffers interface.

a controller comprising hardware logic and/or software logic for performing operations for repeatedly reconfiguring task assignment to the array of processing units, the operations comprising, for each iteration of a plurality of iterations,

302. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Microsoft Azure includes a controller (e.g., at least the AC in conjunction with the RM, SMs, and FMs) comprising hardware logic and/or software logic for

performing operations for repeatedly reconfiguring task assignment to the array of processing units.

303. The controller also repeatedly rearranges communication path connectivity among the array of processing units through the Azure Fabric.

identifying, for each software application of at least a portion of the plurality of software applications, an amount of input data at one or more input data buffers of the plurality of input data buffers buffering data for the respective software application,

304. For each iteration of a plurality of iterations, Microsoft Azure identifies, for each software application of at least a portion of the plurality of software applications, an amount of input data at one or more input data buffers of the plurality of input data buffers buffering data for the respective software application.

305. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Azure identifies data at the data queues that is specific to a respective software application.

allocating, to each software application of the portion of the plurality of software applications, a number of processing units of the array of processing units based at least in part on the amount of input data buffered for the respective software application, and

306. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the data queues act as input buffer, while the RM in coordination with the SMs and FMs, which work in conjunction with the AC, allocates processing units based at least in part on the amount of application program data stored at a queue.

307. Microsoft Azure allocates, to each software application of the portion of the plurality of software applications, a number of processing units of the array of processing units based at least in part on the amount of input data buffered for the respective software application.

308. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the controller (*e.g.*, at least the AC in conjunction with the RM, SMs,

and FMs) allocates accelerators (*e.g.*, FPGA processors) to software applications. Each software application requires processing capacity as reflected by the data load buffered in the queues. The number of processor cores allocated to each application is adjustable based on the application program's processing load. In the illustrative example depicted above, the Ranking Service program is allocated twice the number of cores that are assigned to the ADLA program due to the processing needs of the Ranking Service as compared to the ADLA program.

for each software application of the portion,

- i) assigning one or more task instances of the respective software application for concurrent processing of the amount of input data to the number of processing units allocated to the respective software application by the allocating as one or more assigned instances, and*

309. Microsoft Azure assigns one or more task instances of the respective software application for concurrent processing of the amount of input data to the number of processing units allocated to the respective software application by the allocating as one or more assigned instances.

310. In Azure, one or more task instances is assigned (*e.g.*, through the SM in coordination with the RM and FMs, which work in conjunction with the AC) to different processing units of the one or more processing units such that assigning results in individual processing units of the plurality of processing units being assigned their respective task instances for concurrent processing of the input load for the respective applications.

- ii) adjusting, based at least in part on a change in a count of units between the number of processing units allocated to the respective software application and a number of previously allocated processing units allocated to the respective software application during a previous iteration of the plurality of iterations, a relative portion of the amount of input data to be processed by at least one assigned instance of the one or more assigned instances;*

311. Microsoft Azure adjusts, based at least in part on a change in a count of units between the number of processing units allocated to the respective software application and a number of previously allocated processing units allocated to the respective software application

during a previous iteration of the plurality of iterations, a relative portion of the amount of input data to be processed by at least one assigned instance of the one or more assigned instances.

312. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, a given software application has its own processing load at the input buffers. The number of processor cores allocated to any given application is adjustable based on the application’s processing needs. In the illustrative example depicted above, the Ranking Service program is allocated twice the number of cores than are allocated to the ADLA due to the relative processing needs between the two programs. In that scenario, the number of executing instances assigned per a given program may be adjusted due to the reallocation of cores.

wherein, for one or more iterations of the plurality of iterations where a current number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is increased by the allocating, adjusting comprises relatively decreasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application; and

wherein for one or more other iterations of the plurality of iterations where a present number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is decreased by the allocating, adjusting comprises relatively increasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application.

313. The adjusting performed by Microsoft Azure comprises, for one or more iterations of the plurality of iterations where a current number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is increased by the allocating, *inter alia*, relatively decreasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application. For one or more other iterations of the plurality of iterations where a present number of the one or more processing units allocated to any given software application of the portion of the plurality of software applications is decreased by the allocating, the adjusting

comprises relatively increasing the portion of the amount of input data to be processed by at least one instance of the one or more assigned instances of the respective software application.

314. As discussed above, Azure changes the graphs or groupings of processing units dynamically and periodically (*e.g.*, at least through the AC in conjunction with the RM, SMs, and FMs). In some situations, such as the scenarios described in the '709 patent, in a given reallocation some of the processing units are added to the group of processors performing a task (*e.g.*, mathematical algorithm) for an application (*e.g.*, search ranking). For a given input load level, the amount of input data sent to each processing unit in that group decreases since there are more processor units performing the mathematical algorithm on a given amount of input data to be processed, and vice versa in case fewer processing units were allocated for that given workload.

315. On information and belief, in accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the '242 patent no later than its issuance.

316. Microsoft has had knowledge of and notice of the chain of applications underlying the '242 patent since at least January 2018 when ThroughPuter brought the '090 and '833 patents as well as its portfolio, including a number of pending patent applications, to the attention of Microsoft in the January 2018 letter.

317. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '242 patent.

318. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '242 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into

the United States, the Microsoft Azure platform or components and services thereof.

319. Microsoft's infringement of the '242 patent has been, and continues to be knowing, intentional, and willful. On information and belief, Microsoft has had knowledge of and notice of the application underlying the '242 patent since at least January 2018 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

320. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '242 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '242 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

321. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '242 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

322. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific

intent that its customers will use the Azure platform to infringe the '242 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

323. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the Azure platform.

324. On information and belief, Microsoft has engaged in the above activities with knowledge of the '242 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

325. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '242 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the '242 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '242 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute a material part of the claimed invention of the '242 patent.

326. Microsoft's infringement of the '242 patent has injured ThroughPuter in its

business and property rights. Microsoft's infringement of the '242 patent has been and is deliberate and willful and constitutes egregious misconduct. On information and belief, despite actual knowledge of the '242 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '242 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

327. Microsoft's infringement of the '242 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 5
(Infringement of U.S. Patent No. 10,318,353)

328. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

329. On June 11, 2019, the United States Patent and Trademark Office duly and legally issued the '353 patent, entitled "Concurrent Program Execution Optimization." A copy of the '353 patent is attached as Exhibit 5.

330. Mark Sandstrom is the sole and true inventor of the '353 patent.

331. ThroughPuter owns all right, title, and interest to and in the '353 patent.

332. Microsoft infringes at least claims 1-2 and 8-14 of the '353 patent.

333. Claim 1 of the '353 patent is representative of the claims infringed by Microsoft

and recites:

1. A system for processing a set of computer program instances, comprising:
 a plurality of processing stages, at least one of the plurality of processing stages comprising multiple processing cores, wherein,
 each given task of a plurality of tasks of a given program instance of the set of program instances is hosted at a different stage of the plurality of processing stages as a local task of the given program instance at the respective stage, and
 for at least one of the multiple processing cores of a given processing stage of the plurality of processing stages, a local task of one of the program instances is assigned as an active task instance for execution for a period of time; and
 a group of multiplexers each connecting inter-task communications (ITC) data to a respective stage of the plurality of processing stages, wherein at least one multiplexer of the group of multiplexers is a hardware resource dedicated to the local task, wherein
 the at least one multiplexer is configured to connect ITC data to any processing core of the multiple processing cores to which the local task is assigned for execution for the period of time.

334. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '353 patent.

1. A system for processing a set of computer program instances, comprising:

335. As explained above in the section entitled, "Microsoft's Infringing Cloud Computing Architecture," Microsoft Azure comprises a system for processing a set of computer program instances.

a plurality of processing stages, at least one of the plurality of processing stages comprising multiple processing cores, wherein, each given task of a plurality of tasks of a given program instance of the set of program instances is hosted at a different stage of the plurality of processing stages as a local task of the given program instance at the respective stage, and

336. Microsoft Azure includes a plurality of processing stages where at least one of the plurality of processing stages comprises multiple processing cores, wherein each given task of a plurality of tasks of a given program instance of the set of program instances is hosted at a different stage of the plurality of processing stages as a local task of the given program instance at the

and/or parallelizable, the stage or task may be executed by more than one computing core. This functionality is explained, for example, in Microsoft's '392 patent:

In another scenario, assume that the function has plural stages that the function parsing component 3220 maps to different acceleration components. But assume that one stage is more labor intensive than the others. To avoid a bottleneck in processing associated with this stage, the function parsing component 3220 can allocate two or more acceleration components that operate in parallel for this stage.

Id. at 32:40-47.

for at least one of the multiple processing cores of a given processing stage of the plurality of processing stages, a local task of one of the program instances is assigned as an active task instance for execution for a period of time; and

340. In Microsoft Azure, for at least one of the multiple processing cores of a given processing stage of the plurality of processing stages, a local task of one of the program instances is assigned as an active task instance for execution for a period of time.

341. In Azure, a local task to be executed by a processing stage is assigned as an active task instance for execution on a processing core for a period of time.

a group of multiplexers each connecting inter-task communications (ITC) data to a respective stage of the plurality of processing stages, wherein at least one multiplexer of the group of multiplexers is a hardware resource dedicated to the local task, wherein the at least one multiplexer is configured to connect ITC data to any processing core of the multiple processing cores to which the local task is assigned for execution for the period of time.

342. Microsoft Azure includes a group of multiplexers each connecting inter-task communications (ITC) data to a respective stage of the plurality of processing stages, wherein at least one multiplexer of the group of multiplexers is a hardware resource dedicated to the local task, wherein the at least one multiplexer is configured to connect ITC data to any processing core of the multiple processing cores to which the local task is assigned for execution for the period of time.

343. As explained in the above and further discussed in Microsoft's '392 patent, each of the acceleration components H (2502 in the figure below) communicates with other acceleration components through **multiport switch 2532**, **multiplexers**, and additional switches and **multiplexers** in the top-of-rack (TOR) switch (not shown).

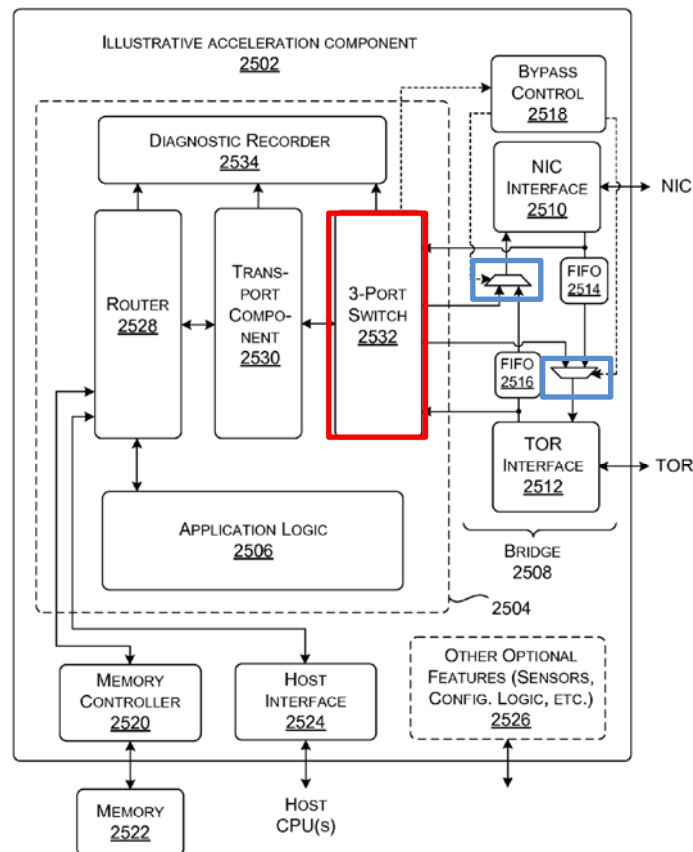


FIG. 25

Exhibit 16 at Fig. 25 (annotated).

344. When one acceleration component needs to send inter-task communications data to another of the pipelined stages executing on a remote acceleration component, the data is sent through **multiport switch 2532** (which is an example of a hardware resource dedicated to the local task), **multiplexers**, and additional switches and **multiplexers** in the TOR switch, through the TOR interface on the remote acceleration component, and to the **multiport switch 2532** on the

remote acceleration component.

345. When one acceleration component needs to receive data from another of the pipelined stages executing on a different acceleration component, the data is sent through **multiport switch 2532** on the remote acceleration component, through the **multiplexers** on the remote acceleration component, through the additional switches and **multiplexers** in the TOR switch, through the TOR interface on 2512 on the local acceleration component, and to the **multiport switch 2532** on the local acceleration component.

346. In accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the '353 patent no later than its issuance.

347. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '353 patent.

348. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '353 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

349. Microsoft's infringement of the '353 patent has been, and continues to be knowing, intentional, and willful. Microsoft has had knowledge of and notice of the application underlying the '353 patent since at least January 2018 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

350. Microsoft actively, knowingly, and intentionally has induced, or has threatened to

induce, infringement of the '353 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '353 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

351. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '353 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

352. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '353 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

353. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the

Azure platform.

354. Microsoft has engaged in the above activities with knowledge of the '353 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

355. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '353 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the '353 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '353 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute a material part of the claimed invention of the '353 patent.

356. Microsoft's infringement of the '353 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '353 patent has been and is deliberate and willful and constitutes egregious misconduct. Despite actual knowledge of the '353 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of

the '353 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

357. Microsoft's infringement of the '353 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 6
(Infringement of U.S. Patent No. 10,310,902)

358. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

359. On June 4, 2019, the United States Patent and Trademark Office duly and legally issued the '902 patent, entitled "System and Method for Input Data Load Adaptive Parallel Processing." A copy of the '902 patent is attached as Exhibit 6.

360. Mark Sandstrom is the sole and true inventor of the '902 patent.

361. ThroughPuter owns all right, title, and interest to and in the '902 patent.

362. Microsoft infringes at least claims 1-2 and 4-16 of the '902 patent.

363. Claim 1 of the '902 patent is representative of the claims infringed by Microsoft and recites:

1. A system, executing on at least one of hardware logic and software logic executing on a plurality of processors, for hosting a plurality of application programs, the system comprising:

a plurality of processing data input buffers, each input buffer of the plurality of processing data input buffers queuing data for a corresponding instance of one program of the plurality of application programs, wherein each program of the plurality of programs comprises a plurality of instances;

an array of cores of computing capacity;

a first subsystem configured to allocate the array of cores of computing capacity among the plurality of application programs, wherein

allocating comprises allocating the array of cores of computing capacity based at least in part on a respective volume of processing data at each buffer of the plurality of processing data input buffers, and a respective processing quota of each application program of at least a portion of the plurality of application programs, and allocating comprises allocating more than one core of the array of cores of computing capacity to at least one of the plurality of application programs;

a second subsystem configured to assign, for each program of the plurality of application programs, each core allocated to the respective program to a different instance of the plurality of instances of the respective program, wherein

the assigning results in assignment of a plurality of selected instances of the plurality of instances of the plurality of application programs, the plurality of instances comprising one or more executable instances of each program of the plurality of application programs, wherein

a number of the plurality of selected instances is fewer than a maximum number of the plurality of instances, and the plurality of selected instances is selected based at least in part on respective volumes of processing data available for each instance of the plurality of instances of the respective program at the portion of the plurality of data input buffers queuing data for the respective program, and according to the assigning, control connectivity between the plurality of processing data input buffers and the array of cores; and

a third subsystem configured, according to the controlling, to establish direct data access from each input buffer of at least a subset of the plurality of processing data input buffers to the respective core of the array of cores that is assigned to a given corresponding instance of the program for which the respective input buffer is queuing data;

wherein the array of cores is periodically allocated by the first subsystem and assigned by the second subsystem based at least in part on changes in respective volumes of processing data associated with each program of the plurality of application programs.

364. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '902 patent.

1. A system, executing on at least one of hardware logic and software logic executing on a plurality of processors, for hosting a plurality of application programs, the system comprising:

365. As explained above in the section entitled "Microsoft's Infringing Cloud

Computing Architecture,” Microsoft Azure comprises a system, executing on at least one of hardware logic and software logic executing on a plurality of processors, for hosting a plurality of application programs. Microsoft Azure is configured to, and Microsoft uses Azure to, manage and host a number of applications across a pool of processing nodes. The AC in conjunction with the RM, SMs, and FMs is used to manage the assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs, and/or FPGA processors.

a plurality of processing data input buffers, each input buffer of the plurality of processing data input buffers queuing data for a corresponding instance of one program of the plurality of application programs, wherein each program of the plurality of programs comprises a plurality of instances;

366. Microsoft Azure includes a plurality of processing data input buffers, each input buffer of the plurality of processing data input buffers queuing data for a corresponding instance of one program of the plurality of application programs, wherein each program of the plurality of programs comprises a plurality of instances.

367. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” in Azure, queues are associated with specific applications, which act as input buffers for queuing data provided by application programs.

an array of cores of computing capacity;

368. Microsoft Azure includes an array of cores of computing capacity.

369. As discussed above, Azure includes an array of processing cores (*e.g.*, FPGA processors).

a first subsystem configured to allocate the array of cores of computing capacity among the plurality of application programs, wherein allocating comprises allocating the array of cores of computing capacity based at least in part on a respective volume of processing data at each buffer of the plurality of processing data input buffers, and a respective processing quota of each application program of at least a portion of the plurality of application

programs, and allocating comprises allocating more than one core of the array of cores of computing capacity to at least one of the plurality of application programs;

370. Microsoft Azure includes a first subsystem configured to allocate the array of cores of computing capacity among the plurality of application programs, wherein allocating comprises allocating the array of cores of computing capacity based at least in part on a respective volume of processing data at each buffer of the plurality of processing data input buffers, and a respective processing quota of each application program of at least a portion of the plurality of application programs, and allocating comprises allocating more than one core of the array of cores of computing capacity to at least one of the plurality of application programs.

371. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the RM in conjunction with the SMs and FMs, which work in coordination with the AC, periodically reallocates processing cores among various application programs (*e.g.*, web search ranking, DNN, SQL, SDN offload) based on the processing needs of each application.

372. Microsoft Azure also supports subscriptions for assured processing resource access levels for an additional fee. The RM in conjunction with the SMs and FMs, which work in coordination with the AC, allocates processing resources based in part on the amount of resources subscribed to by a given client for its applications.

a second subsystem configured to assign, for each program of the plurality of application programs, each core allocated to the respective program to a different instance of the plurality of instances of the respective program, wherein the assigning results in assignment of a plurality of selected instances of the plurality of instances of the plurality of application programs, the plurality of instances comprising one or more executable instances of each program of the plurality of application programs, wherein a number of the plurality of selected instances is fewer than a maximum number of the plurality of instances, and the plurality of selected instances is selected based at least in part on respective volumes of processing data available for each instance of the plurality of instances of the respective program at the portion of the plurality of

data input buffers queuing data for the respective program, and according to the assigning, control connectivity between the plurality of processing data input buffers and the array of cores; and

373. Microsoft Azure includes a second subsystem configured to assign, for each program of the plurality of application programs, each core allocated to the respective program to a different instance of the plurality of instances of the respective program.

374. The assigning results in assignment of a plurality of selected instances of the plurality of instances of the plurality of application programs, the plurality of instances comprising one or more executable instances of each program of the plurality of application programs.

375. A number of the plurality of selected instances is fewer than a maximum number of the plurality of instances, and the plurality of selected instances is selected based at least in part on respective volumes of processing data available for each instance of the plurality of instances of the respective program at the portion of the plurality of data input buffers queuing data for the respective program, and according to the assigning, control connectivity between the plurality of processing data input buffers and the array of cores.

376. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the SMs in coordination with the RM and FMs, which work in conjunction with the AC, select tasks to be assigned to, and executed by, each processing core in a group. The selection and assignment of the respective application instances is a function of the number of cores in the group, the readiness of the application task instances for execution, and the processing needs of each application instance as manifested at the input buffers or queues for each program. A given application will often not have all of its possible tasks processed all the time.

377. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the controller through the Azure Fabric controls connectivity between

the processing cores.

a third subsystem configured, according to the controlling, to establish direct data access from each input buffer of at least a subset of the plurality of processing data input buffers to the respective core of the array of cores that is assigned to a given corresponding instance of the program for which the respective input buffer is queuing data; wherein the array of cores is periodically allocated by the first subsystem and assigned by the second subsystem based at least in part on changes in respective volumes of processing data associated with each program of the plurality of application programs.

378. Microsoft Azure includes a third subsystem configured, according to the controlling, to establish direct data access from each input buffer of at least a subset of the plurality of processing data input buffers to the respective core of the array of cores that is assigned to a given corresponding instance of the program for which the respective input buffer is queuing data.

379. Azure's FMs in coordination with SMs and RM, which work in conjunction with the AC, establish direct data input access from the queues to the appropriate processing cores.

380. In Microsoft Azure, the array of cores is periodically allocated by the first subsystem and assigned by the second subsystem based at least in part on changes in respective volumes of processing data associated with each program of the plurality of application programs.

381. As discussed above in the section entitled "Microsoft's Infringing Cloud Computing Architecture," the RM in coordination with the SMs and FMs (which work in conjunction with the AC) allocates and assigns processors based on changes in processing volume. In some situations, such as the scenarios described in the '709 patent, in a given reallocation some of the processor cores are added to the group of processors performing a task (*e.g.*, mathematical algorithm) for an application (*e.g.*, search). An impact of this is that the amount of input data sent to individual cores in that group decreases since there are more cores performing the mathematical algorithm per a given volume of input data to be processed (*e.g.*, searched for key terms).

382. Upon information and belief, in accordance with 35 U.S.C. § 287, Microsoft has

had actual notice and knowledge of the '902 patent no later than its issuance.

383. Microsoft has had knowledge of and notice of the chain of applications underlying the '902 patent since at least January 2018 when ThroughPuter brought the '090 and '833 patents as well as its portfolio, including a number of pending patent applications, to the attention of Microsoft in the January 2018 letter.

384. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '902 patent.

385. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '902 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

386. Microsoft's infringement of the '902 patent has been, and continues to be knowing, intentional, and willful. On information and belief, Microsoft has had knowledge of and notice of the application underlying the '902 patent since at least January 2018 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

387. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '902 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '902 patent by executing the system

operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

388. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '902 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

389. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '902 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

390. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the Azure platform.

391. On information and belief, Microsoft has engaged in the above activities with knowledge of the '902 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

392. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '902 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the '902 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '902 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute a material part of the claimed invention of the '902 patent.

393. Microsoft's infringement of the '902 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '902 patent has been and is deliberate and willful and constitutes egregious misconduct. On information and belief, despite actual knowledge of the '902 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '902 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

394. Microsoft's infringement of the '902 patent is exceptional and entitles

ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 7
(Infringement of U.S. Patent No. 10,133,599)

395. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

396. On November 20, 2018, the United States Patent and Trademark Office duly and legally issued the '599 patent, entitled "Application Load Adaptive Multi-Stage Parallel Data Processing Architecture." A copy of the '599 patent is attached as Exhibit 7.

397. Mark Sandstrom is the sole and true inventor of the '599 patent.

398. ThroughPuter owns all right, title, and interest to and in the '599 patent.

399. Microsoft infringes at least claims 1-8, 10-11, 13-15, and 17-20 of the '599 patent.

400. Claim 1 of the '599 patent is representative of the claims infringed by Microsoft and recites:

1. A system including processors and comprising a plurality of subsystems, implemented on hardware logic and/or software logic executing on the processors, for dynamic resource management of a pool of processing resources on behalf of application programs, the system comprising:
 - a first subsystem configured to periodically allocate a plurality of processing units of the pool of processing resources among a plurality of application programs over time, wherein
 - the plurality of processing units comprises units of at least two types of processing units and/or units of a reconfigurable type, and
 - allocating the plurality of processing units among the plurality of application programs is based at least in part on i) a respective processing demand of each of the plurality of application programs, and ii) a respective processing resource quota of each program of a subset of the plurality of application programs;
 - a second subsystem configured to, for each program of the plurality of application programs, select a set of highest priority instances of one or more instances of a respective program, wherein

a number of instances in the set of instances corresponds to a number of processing units allocated to the respective application program by the first subsystem during a current allocation period, and each instance of the set of highest priority instances is selected based at least in part on a) the number of processing units allocated to the respective program and b) a relative readiness for execution among the one or more instances of the respective application program; and

a third subsystem configured to assign, for each of the plurality of application programs, the respective set of highest priority instances to a subset of the processing units allocated by the first subsystem for the respective application program for execution during an upcoming allocation period, wherein, for at least a portion of instances of at least a portion of the plurality of application programs,

a respective instance of the respective set of highest priority instances is associated with a processing unit type or configuration, wherein assigning comprises placing the respective instance to a respective processing unit of the plurality of processing units based at least in part on the processing unit type or configuration associated with the respective instance, and prioritizing placement of the respective instance to the particular type or configuration processing unit demanded;

wherein the system is further configured to periodically adjust allocations of the plurality of processing units of the pool of processing resources over time, wherein, for each periodic adjustment, the adjusting is followed by the selecting and the assigning.

401. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '599 patent.

1. A system including processors and comprising a plurality of subsystems, implemented on hardware logic and/or software logic executing on the processors, for dynamic resource management of a pool of processing resources on behalf of application programs, the system comprising:

402. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Microsoft Azure comprises a system including processors and comprising a plurality of subsystems, implemented on hardware logic and/or software logic executing on the processors, for dynamic resource management of a pool of processing resources on behalf of application programs.

403. Microsoft Azure is configured to, and Microsoft uses Azure to, manage and host a

number of applications across a pool of processing nodes. The Azure RM (in coordination with the SMs and FMs) is used to manage the assignment of computing cores across multiple applications and multiple application tasks running, for instance, on multiple computing cores such as CPUs, GPUs, and/or FPGA processors.

a first subsystem configured to periodically allocate a plurality of processing units of the pool of processing resources among a plurality of application programs over time, wherein

the plurality of processing units comprises units of at least two types of processing units and/or units of a reconfigurable type, and

allocating the plurality of processing units among the plurality of application programs is based at least in part on i) a respective processing demand of each of the plurality of application programs, and ii) a respective processing resource quota of each program of a subset of the plurality of application programs;

404. Microsoft Azure includes a first subsystem configured to periodically allocate a plurality of processing units of the pool of processing resources among a plurality of application programs over time, wherein the plurality of processing units comprises units of at least two types of processing units and/or units of a reconfigurable type, and allocating the plurality of processing units among the plurality of application programs is based at least in part on i) a respective processing demand of each of the plurality of application programs, and ii) a respective processing resource quota of each program of a subset of the plurality of application programs.

405. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the RM in coordination with the SMs and FMs, which work in conjunction with the AC, periodically reallocates processing units (*e.g.*, CPUs, reconfigurable FPGA cores, and/or other core types) among various application programs (*e.g.*, web search ranking, DNN, SQL, SDN offload) based on the processing demand of each application.

406. Microsoft Azure also supports subscribing for assured resources such as FPGA

Groupings for an additional fee. In such scenarios, the RM in coordination with the SMs and FMs, which work in conjunction with the AC, allocates processors based in part on the number of FPGA Groupings subscribed by the client for its application.

a second subsystem configured to, for each program of the plurality of application programs, select a set of highest priority instances of one or more instances of a respective program, wherein

a number of instances in the set of instances corresponds to a number of processing units allocated to the respective application program by the first subsystem during a current allocation period, and each instance of the set of highest priority instances is selected based at least in part on a) the number of processing units allocated to the respective program and b) a relative readiness for execution among the one or more instances of the respective application program; and

407. Microsoft Azure includes a second subsystem configured to, for each program of the plurality of application programs, select a set of highest priority instances of one or more instances of a respective program, wherein a number of instances in the set of instances corresponds to a number of processing units allocated to the respective application program by the first subsystem for a current allocation period, and each instance of the set of highest priority instances is selected based at least in part on a) the number of processing units allocated to the respective program and b) a relative readiness for execution among the one or more instances of the respective application program.

408. As discussed above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Azure schedules tasks based on their relative readiness for execution. Azure prioritizes tasks based on, *e.g.*, whether a task is awaiting input data and whether data sufficient to execute the task exists in the queue.

409. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the AC working in conjunction with the RM, SMs, and FMs selects tasks to be assigned to, and executed by, each FPGA or other processor core in a group. The

number and selection of the respective tasks is a function of at least the number of processor cores in the group and the respective readiness of the tasks for execution.

a third subsystem configured to assign, for each of the plurality of application programs, the respective set of highest priority instances to a subset of the processing units allocated by the first subsystem for the respective application program for execution during an upcoming allocation period, wherein, for at least a portion of instances of at least a portion of the plurality of application programs,

a respective instance of the respective set of highest priority instances is associated with a processing unit type or configuration, wherein assigning comprises placing the respective instance to a respective processing unit of the plurality of processing units based at least in part on the processing unit type or configuration associated with the respective instance, and prioritizing placement of the respective instance to the particular type or configuration processing unit demanded;

410. Microsoft Azure includes a third subsystem configured to assign, for each of the plurality of application programs, the respective set of highest priority instances to a subset of the processing units allocated by the first subsystem for the respective application program for execution during an upcoming allocation period, wherein, for at least a portion of instances of at least a portion of the plurality of application programs, a respective instance of the respective set of highest priority instances is associated with a processing unit type or configuration, wherein assigning comprises placing the respective instance to a respective processing unit of the plurality of processing units based at least in part on the processing unit type or configuration associated with the respective instance, and prioritizing placement of the respective instance to the particular type or configuration processing unit demanded.

411. The SMs in coordination with the FMs and RM, which work in conjunction with the AC, assign the selected tasks for execution on the individual FPGA or other processor cores in the group.

412. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the FPGA cores in each group are configured or reconfigured to

optimally perform a certain task such as encryption or performing a mathematical algorithm associated with web search ranking. In the use cases, illustrated in the materials discussed above, every task is assigned to a FPGA core specially configured for the task.

wherein the system is further configured to periodically adjust allocations of the plurality of processing units of the pool of processing resources over time, wherein, for each periodic adjustment, the adjusting is followed by the selecting and the assigning.

413. Microsoft Azure is further configured to periodically adjust allocations of the plurality of processing units of the pool of processing resources over time, wherein, for each periodic adjustment, the adjusting is followed by the selecting and the assigning.

414. As discussed above, the RM in coordination with the SMs and FMs, allocates a number of FPGA cores among the applications and periodically adjusts that allocation while the SMs (which work in conjunction with the AC, RM, and FMs), select tasks for execution by the appropriately configured FPGA cores.

415. Microsoft has had knowledge of and notice of the chain of applications underlying the '599 patent since at least January 2018 when ThroughPuter brought the '090 and '833 patents as well as the portfolio, including a number of pending patent applications, to the attention of Microsoft in the January 2018 letter.

416. Microsoft has had knowledge of and notice of the application that led to the issuance of the '599 patent since September 2018 when ThroughPuter brought the application to the attention of Microsoft in the September 2018 letter.

417. On information and belief, in accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge of the '599 patent no later than its issuance.

418. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '599 patent.

419. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '599 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

420. Microsoft's infringement of the '599 patent has been, and continues to be knowing, intentional, and willful. On information and belief, Microsoft has had knowledge of and notice of the application underlying the '599 patent since at least January 2018 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

421. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '599 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '599 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

422. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '599 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources

on behalf of application programs of customers of the Azure platform cloud computing services.

423. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and “how-to” materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the ’599 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

424. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers’ infringing use of the Azure platform.

425. On information and belief, Microsoft has engaged in the above activities with knowledge of the ’599 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

426. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the ’599 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed invention of the ’599 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft’s customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure

cloud platform, thereby infringing the '599 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute a material part of the claimed invention of the '599 patent.

427. Microsoft's infringement of the '599 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '599 patent has been and is deliberate and willful and constitutes egregious misconduct. On information and belief, despite actual knowledge of the '599 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '599 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

428. Microsoft's infringement of the '599 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 8
(Infringement of U.S. Patent No. 9,632,833)

429. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

430. On April 25, 2017, the United States Patent and Trademark Office duly and legally issued the '833 patent, entitled "Scheduling Application Instances to Processor Cores Over

Consecutive Allocation Periods Based on Application Requirements.” A copy of the ’833 patent is attached as Exhibit 8.

431. Mark Sandstrom is the sole and true inventor of the ’833 patent.

432. ThroughPuter owns all right, title, and interest to and in the ’833 patent.

433. Microsoft infringes at least claims 10, 13-15, 17-19, 21-22, 24-25, 27, 29, 32-34, and 37 of the ’833 patent.

434. Claim 34 of the ’833 patent is representative of the claims infringed by Microsoft and recites:

34. A system for assigning instances of software programs to an array of processor cores comprising:

a first hardware or software logic subsystem configured to, for each of a series of successive core allocation periods (CAPs), select, from a group of executable instances of a set of software programs, a subset of the executable instances, referred to as selected instances, for execution on the cores of the array for an upcoming CAP, wherein the selection of the selected instances is based, at least in part, on a respective capacity demand indication of each of the set of software programs, with said indication of a given program (a) being based at least in part on a number of its executable instances that presently have input data available for processing and (b) indicating processor core types demanded by its executable instances;

a second hardware or software logic subsystem configured to assign each of the selected instances for execution on a processor core in the array of processor cores based, at least in part, on matching the respective demanded processor core types associated with the selected instances with types of processor cores available for assignment; and

a third subsystem, comprising the array of processor cores, configured to execute the selected instances on their assigned cores over the next CAP, at least in part, to process the input data.

435. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 34 of the ’833 patent.

34. A system for assigning instances of software programs to an array of processor cores comprising:

436. As explained above in the section entitled “Microsoft’s Infringing Cloud

Computing Architecture,” Microsoft Azure comprises a system for assigning instances of software programs to an array of processor cores.

437. Microsoft Azure is configured to, and Microsoft uses Azure to assign instances of software programs to an array of processor cores.

a first hardware or software logic subsystem configured to, for each of a series of successive core allocation periods (CAPs), select, from a group of executable instances of a set of software programs, a subset of the executable instances, referred to as selected instances, for execution on the cores of the array for an upcoming CAP, wherein the selection of the selected instances is based, at least in part, on a respective capacity demand indication of each of the set of software programs, with said indication of a given program (a) being based at least in part on a number of its executable instances that presently have input data available for processing and (b) indicating processor core types demanded by its executable instances;

438. Microsoft Azure includes a first hardware or software logic subsystem configured to, for each of a series of successive core allocation periods (CAPs), select, from a group of executable instances of a set of software programs, a subset of the executable instances, referred to as selected instances, for execution on the cores of the array for an upcoming CAP, wherein the selection of the selected instances is based, at least in part, on a respective capacity demand indication of each of the set of software programs, with said indication of a given program (a) being based at least in part on a number of its executable instances that presently have input data available for processing and (b) indicating processor core types demanded by its executable instances.

439. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the AC, which works in conjunction with the RM, SMs, and FMs, selects a subset of available executable instances for execution on the FPGA processors. The executable instances are selected based on demand indications of respective software programs with such indications being based in part on the number of instances of the respective program

ready for execution and the type of processor core needed for a given instance. For example, the demand indication of a speech translation software will depend on how many tasks can be executed concurrently and the type of processor core(s) that are best-suited to execute these tasks.

a second hardware or software logic subsystem configured to assign each of the selected instances for execution on a processor core in the array of processor cores based, at least in part, on matching the respective demanded processor core types associated with the selected instances with types of processor cores available for assignment; and

440. Microsoft Azure includes a second hardware or software logic subsystem configured to assign each of the selected instances for execution on a processor core in the array of processor cores based, at least in part, on matching the respective demanded processor core types associated with the selected instances with types of processor cores available for assignment.

441. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the SMs in coordination with the RM and FMs, which work in conjunction with the AC, assign the selected instances for execution on an FPGA or FPGA Groupings at least in part by matching the selected instance with a processor core suited to performing the task.

a third subsystem, comprising the array of processor cores, configured to execute the selected instances on their assigned cores over the next CAP, at least in part, to process the input data.

442. Microsoft Azure includes a third subsystem, comprising the array of processor cores, configured to execute the selected instances on their assigned cores over the next CAP, at least in part, to process the input data.

443. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” Azure includes the Azure Fabric, which includes an array of FPGA cores that execute selected instances to process data received from the queues.

444. In accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge

of the '833 patent no later than the time when ThroughPuter informed Microsoft of its existence as set forth, *supra*.

445. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '833 patent.

446. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '833 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

447. Microsoft's infringement of the '833 patent has been, and continues to be knowing, intentional, and willful. Microsoft has had knowledge of and notice of the application underlying the '833 patent since at least May 2015 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

448. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '833 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '833 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

449. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '833 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

450. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '833 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

451. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customer's infringing use of the Azure platform.

452. Microsoft has engaged in the above activities with knowledge of the '833 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

453. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '833 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed

invention of the '833 patent, including at least the Azure PaaS (platform-as-a-service). When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '833 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute at least a material part of the claimed invention of the '833 patent.

454. Microsoft's infringement of the '833 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '833 patent has been and is deliberate and willful and constitutes egregious misconduct. Despite actual knowledge of the '833 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '833 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

455. Microsoft's infringement of the '833 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

COUNT 9
(Infringement of U.S. Patent No. 9,424,090)

456. ThroughPuter repeats and realleges each and every allegation contained above as though fully set forth herein.

457. On August 23, 2016, the United States Patent and Trademark Office duly and legally issued the '090 patent, entitled "Scheduling Tasks to Configurable Processing Cores Based on Task Requirements and Specification." A copy of the '090 patent is attached as Exhibit 9.

458. Mark Sandstrom is the sole and true inventor of the '090 patent.

459. ThroughPuter owns all right, title, and interest to and in the '090 patent.

460. Microsoft infringes at least claims 1-7 and 10-15 of the '090 patent.

461. Claim 1 of the '090 patent is representative of the claims infringed by Microsoft and recites:

1. A method for assigning a set of processing tasks to an array of processor cores of configurable types, the method comprising:

executing time variable subsets of the processing tasks of differing types on their assigned processor cores of matching types, wherein the matching type for the assigned processor core for a given processing task of the set corresponds to a type of a processor core demanded by the given processing task;

for each of a series of core allocation periods (CAPs), selecting, from the set of processing tasks, specific tasks, referred to as selected tasks, for execution on the processor cores for a next CAP at least in part based on core capacity demand expressions associated with the processing tasks;

assigning the selected tasks for execution on the processor cores for the next CAP in a manner to maximize, within the array, a number of processor cores whose assigned processing tasks for a present and the next CAP demands the same type of processor core; and

configuring the array such that the type of any given processor core in the array matches the type of processing task assigned for execution on the given processor core for the next CAP.

462. On information and belief, Azure is implemented in a manner that meets each and every limitation of claim 1 of the '090 patent.

1. A method for assigning a set of processing tasks to an array of processor cores of configurable types, the method comprising:

463. As explained above in the section entitled "Microsoft's Infringing Cloud Computing Architecture," Microsoft Azure assigns a set of processing tasks to an array of

processor cores of configurable types.

executing time variable subsets of the processing tasks of differing types on their assigned processor cores of matching types, wherein the matching type for the assigned processor core for a given processing task of the set corresponds to a type of a processor core demanded by the given processing task;

464. Microsoft Azure executes time variable subsets of the processing tasks of differing types on their assigned processor cores of matching types, wherein the matching type for the assigned processor core for a given processing task of the set corresponds to a type of a processor core demanded by the given processing task.

465. As explained above in the section entitled “Microsoft’s Infringing Cloud Computing Architecture,” the FPGA accelerator core blocks in each group are repeatedly configured or reconfigured to optimally perform tasks, of varying duration, such as encryption or performing a mathematical algorithm associated with web search ranking. In the use cases, illustrated in the materials discussed above, a given task is assigned to an FPGA accelerator configured for optimal execution of the given the task.

for each of a series of core allocation periods (CAPs), selecting, from the set of processing tasks, specific tasks, referred to as selected tasks, for execution on the processor cores for a next CAP at least in part based on core capacity demand expressions associated with the processing tasks;

466. For each of a series of core allocation periods (CAPs), Microsoft Azure selects, from the set of processing tasks, specific tasks for execution on the processor cores for a next CAP at least in part based on core capacity demand expressions associated with the processing tasks.

467. As discussed above, in Azure, the number of FPGA processor core blocks allocated to each application changes based on the demand expressions as observed by Azure’s resource management systems. In the illustrative example depicted above, the Ranking Service program is allocated twice the number of cores that are assigned to the ADLA program due to the higher processing needs of the former.

assigning the selected tasks for execution on the processor cores for the next CAP in a manner to maximize, within the array, a number of processor cores whose assigned processing tasks for a present and the next CAP demands the same type of processor core; and

468. Microsoft Azure assigns the selected tasks for execution on the processor cores for the next CAP in a manner to maximize, within the array, a number of processor cores whose assigned processing tasks for a present and the next CAP demands the same type of processor core.

469. As discussed above, the SMs in coordination with the RM and FMs, which work in conjunction with the AC, assigns the selected tasks to FPGA processors for execution, as illustrated in Microsoft's '709 patent, such that the FPGA cores are not reassigned between tasks demanding different core types if doing so is not necessary. In Azure, FPGA core assignments, at least across tasks demanding distinct core types, are held constant from one cycle to the next to the extent permitted by the variations in selections of tasks being assigned for execution.

configuring the array such that the type of any given processor core in the array matches the type of processing task assigned for execution on the given processor core for the next CAP.

470. Microsoft Azure configures the FPGA core array such that the type of any given processor core in the array matches the type of processing task assigned for execution on the given processor core for the next CAP. Azure's SMs in coordination with the FMs and RM, which work in conjunction with the AC, assign the selected tasks for execution on appropriate individual FPGA cores in the group.

471. As discussed above, the FPGA cores in each group are configured or reconfigured to optimally perform a certain task such as encryption or performing a mathematical algorithm associated with web search ranking. In the use cases, illustrated in the materials discussed above, every task is assigned to an FPGA core specially configured for the task.

472. In accordance with 35 U.S.C. § 287, Microsoft has had actual notice and knowledge

of the '090 patent no later than the time when ThroughPuter informed Microsoft of its existence as set forth, *supra*.

473. On information and belief, Microsoft continues without license to make, use, import, market, offer for sale, and/or sell in the United States services or products that infringe the '090 patent.

474. Microsoft has directly and indirectly infringed and continues to directly and indirectly infringe the '090 patent by engaging in acts constituting infringement under 35 U.S.C. § 271(a), (b), and/or (c), including but not necessarily limited to one or more of making, using, selling and offering to sell, in this District and elsewhere in the United States, and importing into the United States, the Microsoft Azure platform or components and services thereof.

475. Microsoft's infringement of the '090 patent has been, and continues to be knowing, intentional, and willful. Microsoft has had knowledge of and notice of the application underlying the '090 patent since at least May 2015 and despite this knowledge continues to commit the aforementioned infringing acts. For at least the reasons stated in this paragraph and above, this infringement has been willful.

476. Microsoft actively, knowingly, and intentionally has induced, or has threatened to induce, infringement of the '090 patent through a range of activities. First, on information and belief, Microsoft has induced infringement by controlling the design and development of, offering for sale, and selling the services of the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '090 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

477. Second, on information and belief, Microsoft has induced infringement by its customers through the dissemination of promotional, marketing, and tutorial materials relating to the Azure platform with the knowledge and specific intent that its customers will use the Azure platform to infringe the '090 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

478. Third, on information and belief, Microsoft has induced infringement by its customers through the creation and online posting of tutorial and "how-to" materials for the Azure platform and/or its individual components in the United States with the knowledge and specific intent that its customers will use the Azure platform to infringe the '090 patent by executing the system operations and utilizing the system components to perform dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure platform cloud computing services.

479. Fourth, on information and belief, Microsoft has induced infringement through the distribution of other instructional materials, product manuals, and technical materials with the knowledge and the specific intent to encourage and facilitate its customers' infringing use of the Azure platform.

480. Microsoft has engaged in the above activities with knowledge of the '090 patent and with the specific intent to encourage and cause infringement by its customers, as shown by the allegations set forth above.

481. Microsoft has contributed to, or has threatened to contribute to, the infringement by its customers of the '090 patent by, without authority, selling and offering to sell within the United States cloud computing services and customer support services for practicing the claimed

invention of the '090 patent, including at least the Azure platform as a whole and/or the individual components of the Azure platform. When, for example, the Azure platform is used by Microsoft's customers for the various cloud computing services Microsoft offers, the Azure system operations and system components are used to perform the claimed dynamic resource management of the pool of Azure processing resources on behalf of application programs of customers of the Azure cloud platform, thereby infringing the '090 patent. The Azure platform and/or its individual components, supplied by Microsoft, constitute at least a material part of the claimed invention of the '090 patent.

482. Microsoft's infringement of the '090 patent has injured ThroughPuter in its business and property rights. Microsoft's infringement of the '090 patent has been and is deliberate and willful and constitutes egregious misconduct. Despite actual knowledge of the '090 patent and numerous related patents and applications since at least January 2018, Microsoft continued to develop and offer its infringing products and services. In developing and offering its products and services, Microsoft has been willfully blind to this ongoing infringement. ThroughPuter is entitled to recover monetary damages for the injuries arising from Microsoft's willful infringement pursuant to 35 U.S.C. § 284 in an amount to be determined at trial. Microsoft's infringement of the '090 patent has caused irreparable harm to ThroughPuter and will continue to cause such harm unless and until Microsoft's infringing activities are enjoined by this Court.

483. Microsoft's infringement of the '090 patent is exceptional and entitles ThroughPuter to attorneys' fees and costs incurred in prosecuting this action under 35 U.S.C. § 285.

PRAYER FOR RELIEF

WHEREFORE, ThroughPuter respectfully requests that the Court enter judgment against Microsoft as follows:

- A. An adjudication that Microsoft has infringed one or more claims of the '306, '599, '902, '242, '644, '998, '353, '090, and '833 patents;
- B. An order permanently enjoining Microsoft from further infringement of the '306, '599, '902, '242, '644, '998, '353, '090, and '833 patents;
- C. An award of damages pursuant to 35 U.S.C. § 284;
- D. An order that the damages award be increased up to three times the actual amount assessed, pursuant to 35 U.S.C. § 284;
- E. An award to ThroughPuter of its costs, pre- and post-judgment interest, and reasonable expenses to the fullest extent permitted by law;
- F. A declaration that this case is exceptional pursuant to 35 U.S.C. § 285, and an award of attorneys' fees and costs; and
- G. An award to ThroughPuter of such other and further relief as this Court deems just and proper.

DEMAND FOR JURY TRIAL

Pursuant to Rule 38(b) of the Federal Rules of Civil Procedure, ThroughPuter hereby demands a trial by jury on all issues so triable.

Dated: March 31, 2021

/s/ Wayne M. Helge

Wayne Helge, Esq.
Virginia Bar #71074
Davidson Berquist Jackson & Gowdey
8300 Greensboro Drive
Suite 500
McLean, VA 22102
Telephone: (571) 765-7700
Facsimile: (571) 765-7200
whelge@davidsonberquist.com

Attorney for Plaintiff ThroughPuter, Inc.